

# Statistical Computing

## Hidden Markov Models for Bioinformatics

- Part I -

Uwe Menzel, 2011

[uwe.menzel@matstat.org](mailto:uwe.menzel@matstat.org)

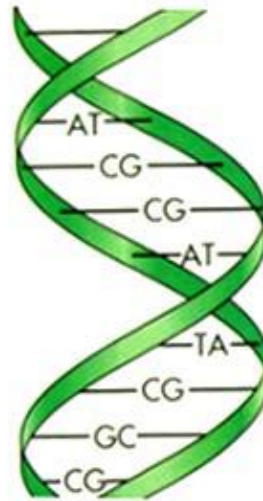
[www.matstat.org](http://www.matstat.org)

# Contents Part I

- What is a Markov chain and what has it to do with DNA?
- A Likelihood Ratio Test using Markov chains to determine whether a small piece of DNA is a CpG island or not
- The Hidden Markov Model: transition and emission probabilities
- Decoding: the Viterbi algorithm
- Forward algorithm, backward algorithm and posterior probabilities
- Parameter estimation for Hidden Markov Models
- A Continuous Density Hidden Markov Model for the recognition of large amplifications and deletions in genomic DNA
- Appendix

# DNA-Sequence

The sequence of the bases (A,T,G,C) in the DNA-molecule determines the blueprint of an organism.

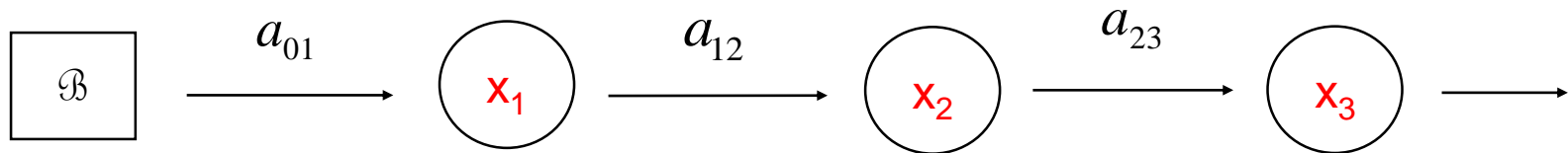


# What is a Markov chain and what has it to do with DNA?



Андрей Андреевич Марков (1856 – 1922)

# Markov chain



## Model:

- a sequence is generated by a random process

## Alphabet:

- set of characters  $x_i$  building up the chain, e.g.  $x_i = \{A, C, G, T\}$

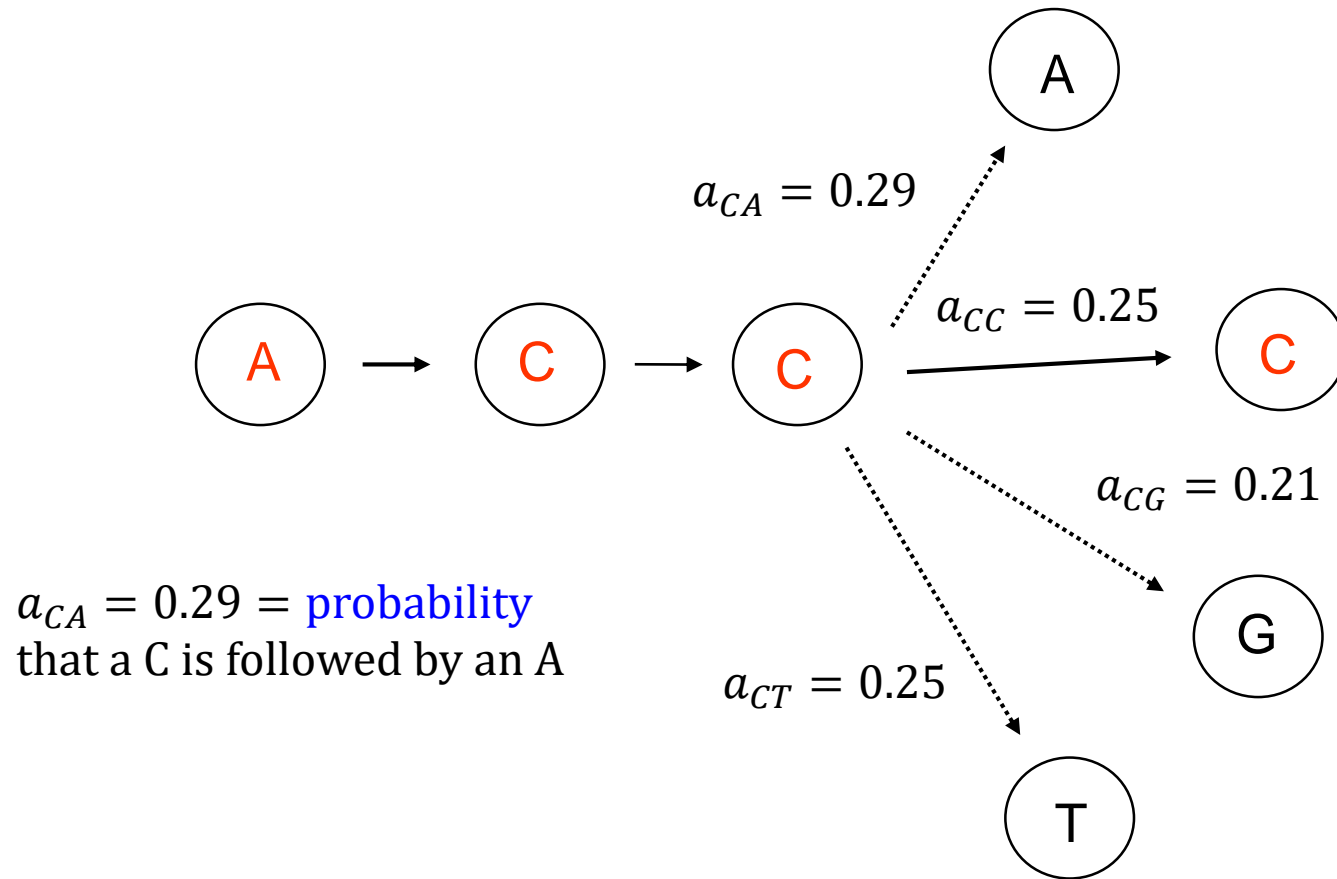
## Markov property:

- the value  $x_{i+1}$  **only** depends on  $x_i$ , but not on  $x_{i-1}, x_{i-2}, \dots$

## Transition probability:

- $a_{st} = P(x_i = t \mid x_{i-1} = s)$

# Markov-chain for DNA



# Drawing a long Markov chain (for DNA)

*B* = begin state

*E* = end state

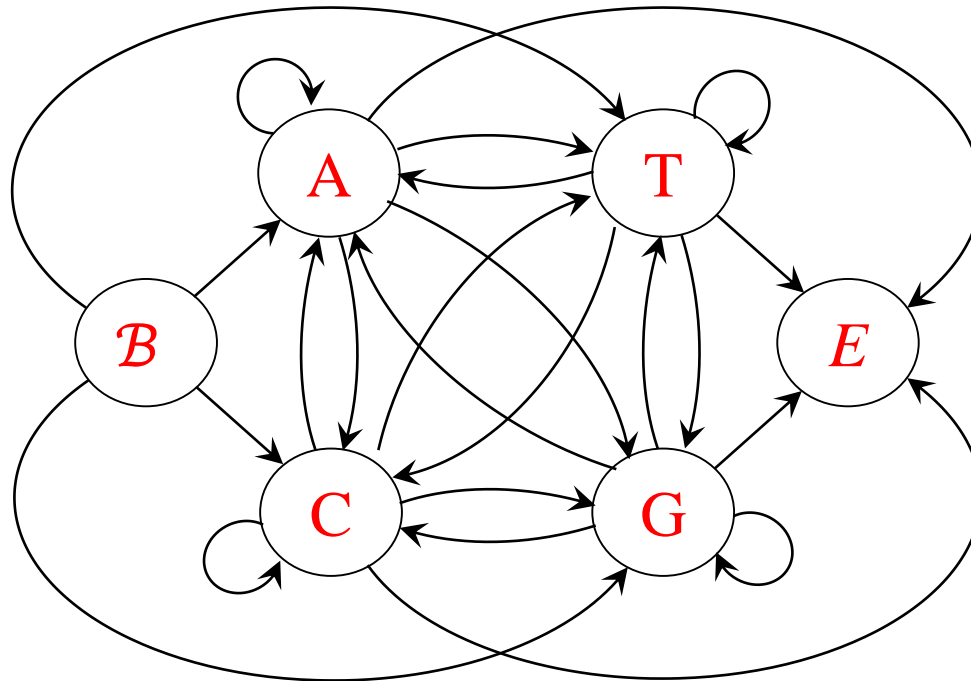


Abbildung: Sven Schuirer

# Joint probability of a chain

$$\begin{aligned} P(\vec{x}) &= P(X_1 = x_1, X_2 = x_2, X_3 = x_3, \dots, X_N = x_N) \\ &= P(x_1, x_2, x_3, \dots, x_N) \end{aligned}$$

$$\begin{aligned} P(\vec{x}) &= P(X_1 = A, X_2 = C, X_3 = C, X_4 = G, X_5 = T) \text{ for DNA} \\ &= P(A, C, C, G, T) \end{aligned}$$

$$P(\vec{x}) = P(x_1, x_2, x_3, \dots, x_{N-1}, x_N)$$

Use multiple times:  $P(x, y) = P(x | y) \cdot P(y)$

$$P(\vec{x}) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_2, x_1) \cdot P(x_4 | x_3, x_2, x_1) \cdot \dots$$

With the **Markov property**, this simplifies to:

$$P(\vec{x}) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_2) \cdot P(x_4 | x_3) \cdot \dots \cdot P(x_N | x_{n-1})$$



# Probability of the Markov chain<sup>1</sup>

$$P(\vec{x}) = P(x_1) \cdot P(x_2 | x_1) \cdot P(x_3 | x_2) \cdot P(x_4 | x_3) \cdot \dots \cdot P(x_N | x_{n-1})$$

Let  $a_{x_{i-1},x_i} = P(x_i | x_{i-1})$  transition probability

$$P(\vec{x}) = P(x_1) \cdot a_{x_1,x_2} \cdot a_{x_2,x_3} \cdot a_{x_3,x_4} \cdot \dots \cdot a_{x_{N-1},x_N}$$

$$P(\vec{x}) = P(x_1) \cdot \prod_{i=2}^N a_{x_{i-1},x_i} \quad \text{with } P(x_1) = a_{x_0,x_1}$$

$x_0$  is a virtual (begin) state, introduced to make the formula nicer

$$P(\vec{x}) = \prod_{i=1}^N a_{x_{i-1},x_i}$$

<sup>1</sup>considering homogeneous Markov-chains only

## Maximum Likelihood (ML) estimators for the transition probabilities in DNA

- count the frequency of dinucleotides,  $c_{st}$ , in many genomic sequences
- normalize: divide by the sum of all outgoing transition probabilities
- a = transition probabilities; c = transition frequencies "counts")

$$a_{st} = \frac{c_{st}}{\sum_i c_{si}} \quad s, t \in \{A, C, G, T\}$$

It follows that:  $\sum_t a_{st} = 1$

The outgoing transition probabilities from each symbol sum up to 1.

For example, it holds that:

$$a_{CA} + a_{CC} + a_{CG} + a_{CT} = 1$$

## Maximum Likelihood (ML) estimators for the transition probabilities in DNA

$$a_{st} = \frac{c_{st}}{\sum_i c_{si}} \quad s, t \in \{A, C, G, T\}$$

$$c_{CG} = 100 \quad c_{CA} = 150 \quad c_{CT} = 50 \quad c_{CC} = 100$$

$$a_{CG} = \frac{c_{CG}}{c_{CG} + c_{CA} + c_{CT} + c_{CC}} = \frac{100}{100 + 150 + 50 + 100} = 0,25$$

$$a_{CA} = \frac{150}{400} = 0,375$$

$$a_{CT} = \frac{50}{400} = 0,125$$

$$a_{CC} = \frac{100}{400} = 0,25$$

$$a_{CG} + a_{CA} + a_{CT} + a_{CC} = 1 \quad \text{line total}$$

# Matrix of transition probabilities

	A	C	G	T
A	0,300	0,205	0,285	0,210
C	0,322	0,298	0,078	0,302
G	0,248	0,246	0,298	0,208
T	0,177	0,239	0,292	0,292

Stochastisc  
matrix

row sum = 1

$$\begin{aligned}P(C, A, A, G) &= a_{0C} \cdot a_{CA} \cdot a_{AA} \cdot a_{AG} \\ &= 0.25 \cdot 0.322 \cdot 0.3 \cdot 0.285 = \underline{\underline{0.00688}}\end{aligned}$$

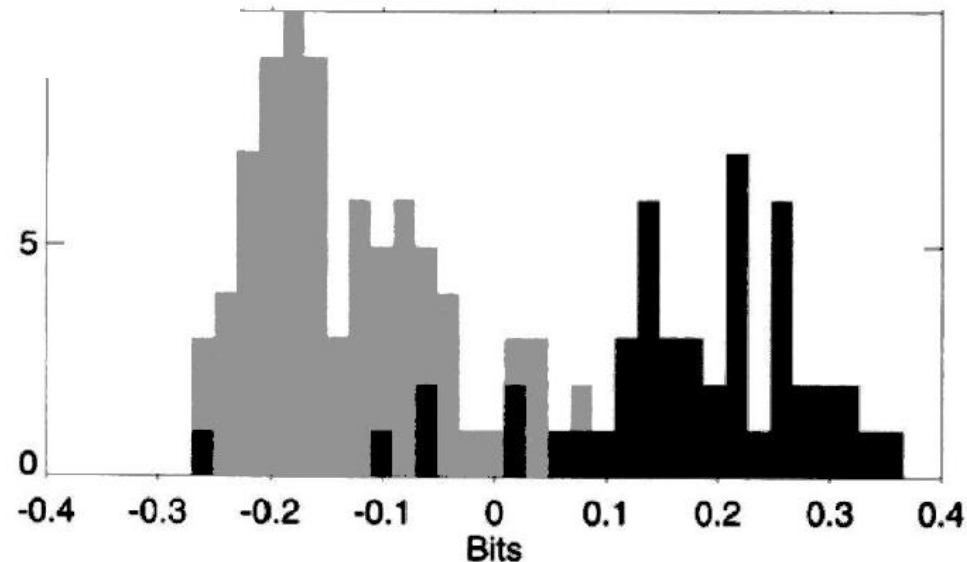
Probability of the chain C, A, A, G ; presupposed the transition probabilities presented in the table, and that the first nucleotide is a C with probability 0.25

# Language recognition (Markov)



- languages differ by the frequencies of transitions between characters:
- for example, "th" is quite frequent in English but not in Spanish

# A Likelihood Ratio Test using Markov chains to determine whether a small piece of DNA is a CpG island or not



Picture from: **Durbin** et al. (Ed): Biological Sequence Analysis, Cambridge University Press, 1998

# What is a CpG island ?



- What CpG frequency do we (approximately) expect ?
  - $P_{CG} \approx \frac{1}{4} \cdot \frac{1}{4} = 1/16$  ; more precisely  $0,21 \cdot 0,21 \approx 4,4\%$
- actual frequency is only **0,8 %** (mammalia)
- cytosine (C) in a CpG is chemically unstable:
  - methylation, deamination:  $CG \rightarrow C^{\text{meth}}G \rightarrow TG$
- **CpG-islands** have a higher CpG percentage, compared to the rest of the genome

**Exon 1 CpG Island: 12634..12767**

```

11941 ttataagatc cccctccctc taaatcctgt ccttctatca cttcatcctt CGctctcctt
12001 taaaatgaga cagttgtcag caggaatcct gCGcaagaac acaccaccct gtttcataga
12061 agatatctca ggtaatgtgc aaacaCGggt ttttaaaCGg agCGcatttt totcatttgt
12121 taatatcacc acctaaatca totcttgcct aaaacaagga gtagaaagtg aatgaaggaa
12181 ggaacaggtg atggtcagtg tctttctac gctcaaaaat ttaagagttt atgtgaaaat
12241 tcataaatat taatctcaat ccagggttaag caaaattttt tgcctcctc tttagaaatt
12301 tctggttgcc aaagttccag aaattgcttc ctcatcctg agcctttcat tttctCGatt
12361 totccattat gtaaCGggga gotggagctt tgggcCGaat ttccaattaa agatgatttt
12421 tacagtcaat gagccaCGtc aggggagCGat ggcaccCGca ggCGgtatca actgatgcaa
12481 gtgttcaagc gaatotoaac tCGtttttct CGgtgactca ttccCGgccc tgcttggcag
12541 CGctgcaccc ttttaacttaa acctCGgcCG gcCGccCGcc gggggcacag agtgtgCGcc
*12601 gggcCGCGCG gcaattggtc ccCGCGcCGa cctcCGccCG CGagCGcCGc CGcttccctt
*12661 cccCGcccCG CGtccctccc cctCGgcccc gCGCGtCGcc tgcctcCGa gccagtCGct
*12721 gacagcCGCG gCGcCGCGag cttctcctct cctcaCGacc gaggcaggta aaCGccCGgg
12781 gtgggaggaa CGCGggCGgg ggcaggggag cCGCGggggc CGagtgagga cccCGggcct
12841 CGggtcccag gCGcaagggt gccCGgcCGg gCGgggtCGg gaccccagtg aggaggggoc
12901 gggggctgcc cCGCGggCGc gtgaCGgtct CGggcctgcc CGgetgCGct ggtctcCGct
12961 CGggtgagge ggcttggctt CGcttttcag gttaggaaag ctccccttac tgCGCGttgg
13021 ggggctgggg gagctggCGg agccaCGtta gggaggtCGg tggCGcCGgg gtgtctcagc
13081 gccccctgca cccCGCGCGg gtcCGgccc aCGggCGatc gctggCGccc agggaaactc
13141 gggagggcCG ccagCGggct cCGcaggCGc ggggCGggga ggggCGcctg ggggcCGCG
13201 ggctCGCGct cccCGccCGt tggcCGcccc tCGgaggcCG agatCGgggc ccagaaCGcc
13261 ccttggcaaa gcttggCGct tcCGCGatgc ccagaggggtg cttgggggga tggagagagg
13321 ggCGccCGcc ggggtagttc CGggagcctc ggtgcctccc gcCGcagctg cagCGttcct
13381 ccCGggagge ggcccagccc ttcctcctCG cCGcctgagc ttctcCGagg ggggctgcag
13441 ccttgCGgcc gttgcccacCG cctggagaag CGgcccCGc ggactgaCGg gCGggggCGg
13501 ggctCGggc ctCGgCGggg gCGgggtcCG gggaggcccc accctctgtt ctccaggggc
13561 ggggagagag gagctgcagg tctgCGgcct ggcccaggt gCGatggCGg accccagctt
13621 ggccagtcac attcctccc gtccccctgg agggagaaCG ctggccatgg ggggctccaa
13681 ggaacaacca gctCGgatg aCGacccttg ggtcacCGgt ctccccacct gtgCGgcagg
13741 CGccttcaCG tttcattatt aaacaatggg gagaaatcca tgtttactgt cctttttagg
13801 aattttttgc totttotott gaggtggctg taggaaatag atttttttt taacctCGca
13861 attcaccac ggtcacatcc atctCGcca tCGcagagcc acagctctcc gtttttgtt
13921 cctagcctcc agattctcac acaacacagt gcagtttcc tgotgtaatg atgaggatc
13981 coatggcCGc gttattttct tgttctgaga gcatcaCGgt ttaattagca gttcccata
14041 tgatttgaag tgtttccCGt ttccttaggg aaaactcctg gtagaatagg attaaggatt
14101 tttacaaata taattatcaa aaacatagga acaggaatt ggataaatat gttaaactc
14161 tggaaaaatc aacaaCGctc ttagatttgt agaagaaagg aaaaaatcac cagtggaaag
14221 gagcaatttt acttacaaa acacagagaa ggtcttacag tgaaaaaaag ctaaccagta

```



# "Training": Finding transition probabilities for both CpG islands and non-islands

## CpG-Islands

	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	<b>0.274</b>	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182

model+

$$a_{CG}^+ = \mathbf{0.274}$$

## Non-Islands

	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	<b>0.078</b>	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

model-

$$a_{CG}^- = \mathbf{0.078}$$

# How to discriminate between CpG islands and non-islands

Observed sequence:  $X = (\text{ATCGCGCGGC})$

$$\begin{aligned} \underline{P(X \mid \text{model}+)} &= \prod_i a_{x_{i-1}x_i}^+ = a_{0A}^+ \cdot a_{AT}^+ \cdot a_{TC}^+ \cdot a_{CG}^+ \cdot a_{GC}^+ \cdot a_{CG}^+ \cdot a_{GC}^+ \cdot a_{CG}^+ \cdot a_{GG}^+ \cdot a_{GC}^+ \\ &= 0.25 \cdot 0.12 \cdot 0.355 \cdot 0.274 \cdot 0.339 \cdot 0.274 \cdot 0.339 \cdot 0.274 \cdot 0.375 \cdot 0.339 \\ &= \underline{3.125 \cdot 10^{-6}} \quad \text{Probability of the chain under CpG island model} \end{aligned}$$

$$\begin{aligned} \underline{P(X \mid \text{model}-)} &= \prod_i a_{x_{i-1}x_i}^- = a_{0A}^- \cdot a_{AT}^- \cdot a_{TC}^- \cdot a_{CG}^- \cdot a_{GC}^- \cdot a_{CG}^- \cdot a_{GC}^- \cdot a_{CG}^- \cdot a_{GG}^- \cdot a_{GC}^- \\ &= 0.25 \cdot 0.21 \cdot 0.239 \cdot 0.078 \cdot 0.246 \cdot 0.078 \cdot 0.246 \cdot 0.078 \cdot 0.298 \cdot 0.246 \\ &= \underline{2.65 \cdot 10^{-8}} \quad \text{Probability of the chain under non-island model} \end{aligned}$$

**Result:** The probability of the chain is higher if we assume model "+" (CpG island model) → **It is more likely that the sequence X is a CpG-Island.**

**Note** that it might be better to carry out these calculations in log-space, to avoid underflow in computations. Products become sums in log-space making the calculations faster.


# Likelihood Ratio Test for Discrimination of CpG islands and non-islands

According to the model, a sequence  $X$  is a CpG-island if:

$$P(X | mod +) > P(X | mod -)$$

$$\frac{P(X | mod +)}{P(X | mod -)} > 1$$

log likelihood ratios

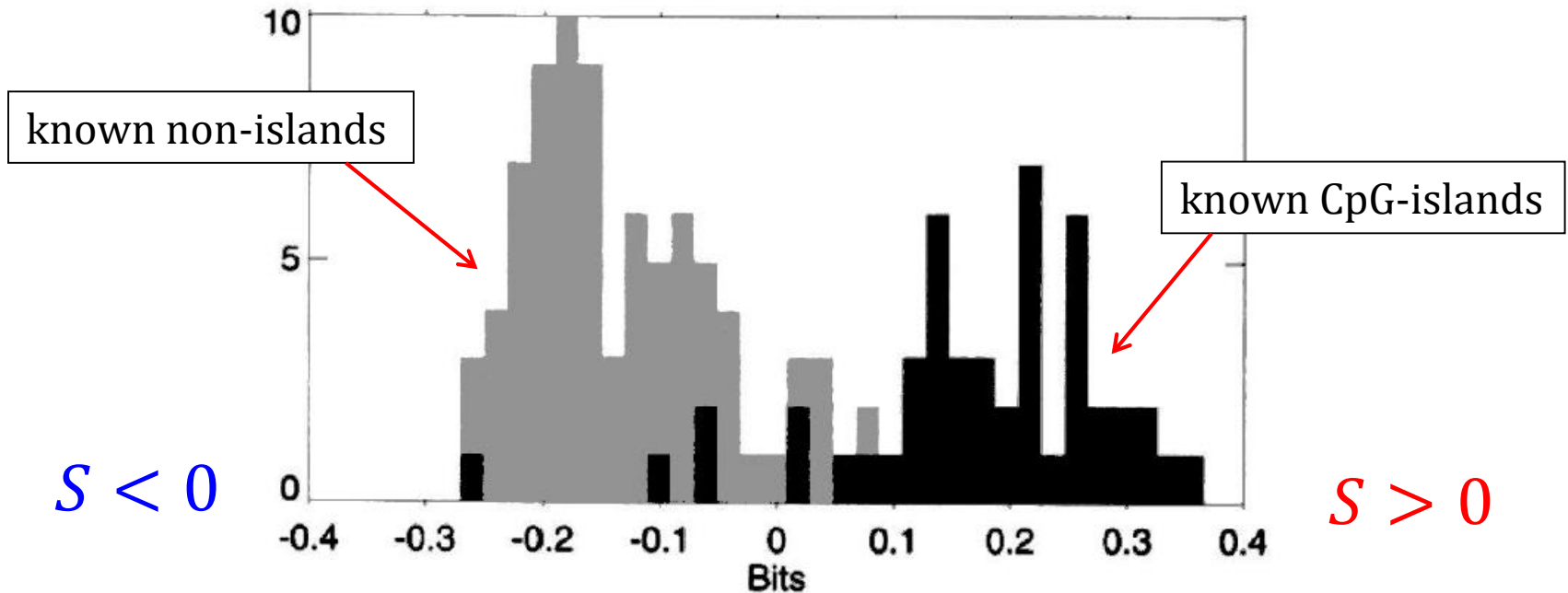
$$S = \log \left[ \frac{P(X | mod +)}{P(X | mod -)} \right] = \sum_{i=1}^L \log \frac{a_{x_{i-1}x_i}^+}{a_{x_{i-1}x_i}^-} = \sum_{i=1}^L \beta_{x_{i-1}x_i} > 0$$


$S =$  log odds score

If the log odds score  $S$  is bigger than 0, it is more likely that the sequence probed is a CpG island. If the log odds score is negative, it is more likely that the sequence does not emerge from a CpG island.

# Does the method really work?

Calculated scores for many training sets (islands and non-islands):



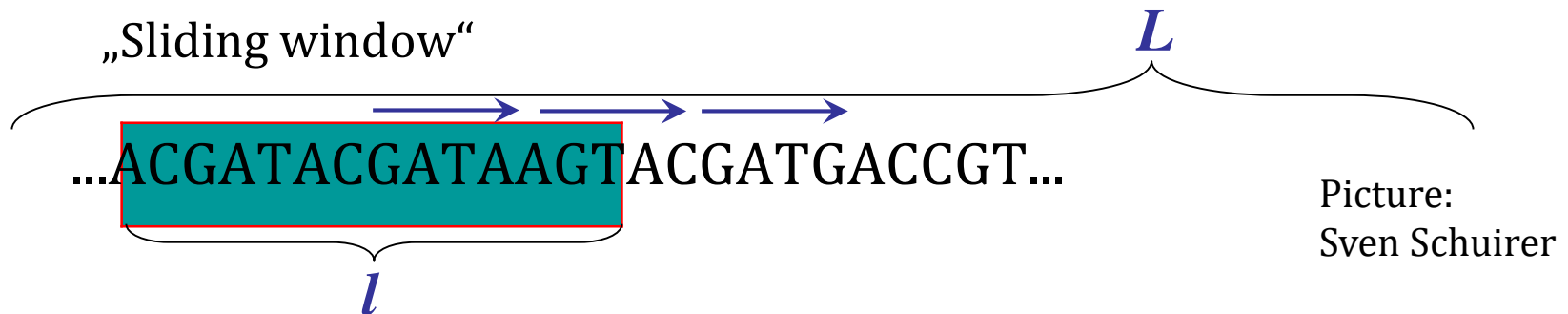
- We see that most of the known CpG islands have a **positive log odds score**, while most of the non-islands have a **negative log odds score**
- **Errors** are caused by: incorrect labels in the training sets, and problems when determining borders between CpG-islands and non-islands

Picture from: Durbin et al. (Ed): Biological Sequence Analysis, Cambridge University Press, 1998

# Pros and Cons of the scoring model

- Given a short piece of DNA, with sufficient certainty, one can decide if it is a CpG-island or not
- You cannot identify a potential CpG-island embedded in a long genomic sequence
- The latter problem can be resolved by using a **Hidden Markov Modell** →

# Long sequence: Finding CpG-islands with a sliding window approach



- Calculate log odds score  $S$  in every window of width  $l$
- Disadvantages:
  - Runtime (?)
  - unknown size of the island  $\rightarrow$  unknown window width  $l$