# **BIOINFORMATICS 1**

or why biologists need computers

 $\underline{http://www.bioinformatics.uni-muenster.de/teaching/courses-2015/bioinf1/index.hbioinf1/index$ 

Prof. Dr. Wojciech Makałowski Institute of Bioinformatics

1







# GENOME ANNOTATION

What are we looking for?

- ·⊱ protein coding genes
- · ≽ RNA coding genes
- gene promoters
- ·⊱ repetitive elements





- Molecular techniques
  - Very laborious
  - Time consuming
  - Expensive
  - Low rate of false positives
- Computational methods
  - Fast
  - Relatively low cost
  - High rate of false positives
  - Poor performance on less typical genes



# GENERAL MODEL OF A GENE

8









12 bioinfo1\_4\_2015 - November 11, 2015



#### MODEL BASED **METHODS**

We take advantage of what we already learned about gene structures and features of coding sequences. Based on this knowledge we can build theoretical model, develop an algorithm to search for important features, train it on known data and use to search for coding sequences in anonymous genomic fragments.

However, we should remember that all models are wrong and only some are useful.



14



15

CODON USAGE

Codon preference in E codiand S typhimuriumgenes.

UCC

UCA

UCG

CCA

ACC ACA ACG

19 17

11 12

Ser Ser

Ser Ser

Pro

Thr

Thr Thr

Thr

UAU

UAG

CAC His

CAA

AAA Lys

AAG Lys

GAU GAC Asp Asp

10

6 8 TTAA STO

9 AAU Ası 17

25 AAC Ası

6 15

16 25

Tyr Tyr

STOP

Gln

12 UGC

13

24 36 12

33 22

IIGA STOP

UGG

CGA Arg 3

AGC Ser 16

AGA Arg

AGG Arg

Cys

Trp

Ser

Gly Gly 24 33 A G U C

Ă

Ċ A G

U C

12

U UUU UUC U

Phe Phe

Leu

Leu

Ile 27 ACU

Ile Ile 27

UUA Leu

UUG

CUU Leu 10 CCU Pro Pro 7 5 CAU His 11 10 CGU Arg Arg 23 23

CUC Leu 10 coc

CUA

CUG Leu 55 COG Pro 15 CAG Gln 31

AUC

AUA

AUG Met 26

GUA Val 12 GCA Ala 16 GA Gl 43 GGI Gl

С

A AUU

G GUU GUC Val Val 17 16 GCU GCC Ala Ala

# SEQUENCE FEATURES

We can check if sequence in particular ORF has some other features which could tell us if this is a putative coding sequence or the ORF is false positive. We can look at the sequence content and compare it with known coding sequence and noncoding sequence and check to which of these two the ORF sequence is more similar to.

#### HIDDEN MARKOV MODELS

- HHM is a statistical model for an ordered sequence of symbols, acting as a stochastic state machine that generates a symbol each time a transition is made from one state to the next. Transitions between states are specified by transition probabilities. A Markov process is a process that moves from state to state depending on the previous n states.
- HHM has been previously used very successfully for speech recognition.
- In biology is used to produce multiple sequence alignments, in generating sequence profiles, to analyze sequence composition and patterns, to produce a protein structure prediction, and to locate genes.
- In gene identification HMM is a model of periodic patterns in a sequence, representing, for example, patterns found in the exons of a gene. HMM provides a measure of how close the data pattern in the sequence resemble the data used to train the model.

19

### MARKOV CHAINS

A Markov Chain is a non-deterministic system in which it is assumed that the probability of moving from one state to another doesn't vary with time. This means the current state and transition does not depend on what happened in the past. The Markov Chain is defined by probabilities for each occurring transition.



20



HOW FAR CAN WE GO?

- Order of our model will have influence on specificity and sensitivity of our program.
  - \* Too short sequences may not be specific enough and program may return a lot of false positives.
  - Long chains may be too specific and our program will not be sensitive enough returning false negatives.

#### ORDER OF MARKOV CHAINS

GCGCTAGCGCCGATCATCTACTCG

GCGCT**AG**CGCCGATCATCTACTCG GCGCTA**GC**GCCGATCATCTACTCG

GCGC**TAG**CGCCGATCATCTACTCG GCGCT**AGC**GCCGATCATCTACTCG

GCGCTAGCGCCGATCATCTACTCG GCGCTAGCGCCGATCATCTACTCG



First order

Second order

Fifth order

For non-coding sequence we assume that probability of each transition is equal. The more 'popular' in coding sequence transition, the higher probability the sequence is coding

#### Probability matrix

Number of probabilities in a DNA matrix of a given order can be calculated according to the following formula:

#### 4<sup>k+1</sup>

where 4 represents number of letters in the DNA alphabet and k stands for the order number.

Hence, first order Markov Model matrix consists of  $4^2 = 16$  probabilities

p(A/A), p(A/T), p(A/C), p(A/G), p(T/A), p(T/T), p(T/C), p(T/G), p(C/A), p(C/T), p(C/C), p(C/G), p(G/A), p(G/T), p(G/C), p(G/G)

25

GCG CTA GCG CCG ATC ATC TAC TCG

G CGC TAG CGC CGA TCA TCT ACT CG

GC GCT AGC GCC GAT CAT CTA CTC G

Frequencies of transitions may depend on in which codon position (1st, 2nd, or 3rd) is a given nucleotide (state)

26

#### Number of probabilities

Codon position 1	Codon position 2	Codon position 3
ACGT	ACGT	A C G T
A .36 .27 .35 .18	A .16 .19 .15 .07	A .22 .33 .24 .13
C .21 .23 .24 .27	C .28 .44 .41 .33	C .21 .29 .27 .21
G .19 .14 .23 .23	G .40 .12 .27 .45	G .44 .15 .37 .53
T .24 .35 .19 .31	T .16 .25 .17 .16	T .13 .22 .12 .13

$$4^{1+1} = 4^2 = 16$$
  
3 (4<sup>1+1</sup>)= 3 x 4<sup>2</sup> = 48

27

#### CALCULATING CODING POTENTIAL OF A GIVEN SEQUENCE

To estimate if the sequence is coding we have to calculate probability that sequence is coding and probability the sequence is non-coding. Next we calculate logarithm from the ratio of these two probability values.

$$LP(S) = \log \frac{Pi(S)}{P_0(S)}$$

If the calculated value is > 0 the likelihood that the sequence is coding is higher than the sequence is not coding, if value is < 0 there is higher likelihood that sequence is not coding.

28

# CODING VS. NON CODING SEQUENCE

A/A	C/A	G/A	T/A coding
0.36	0.21	0.19	0.24
A/A	C/A	G/A	T/A non coding
0.25	0.25	0.25	0.25

$LP(S) = \log \frac{P'(S)}{P_0(S)}$	Codon position 1 A C G T A .36 .27 .35 .18 C 21 .23 .24 .27	Codon position 2 A C G T A .16 .19 .15 .07 C 28 .44 .41 .33	Codon position 3 A C G T A .22 .33 .24 .13 C 21 .29 .27 .21		
S=AGGACG	G .19 .14 .23 .23 T .24 .35 .19 .31	G .40 .12 .27 .45 T .16 .25 .17 .16	G .44 .15 .37 .53 T .13 .22 .12 .13		
$P(S)^{1} = f(A,1)^{*}F(G,A)F(G,G)F(A,G)F(C,A)F(G,C)$					
$P(S) = 0.27 \ge 0.19 \ge 0.27$	x 0.24 x 0.21	x 0.12 = 0.000	008377		
P(S) = 0.25 x 0.25 x 0.25 x 0.25 x 0.25 x 0.25 = 0.0002441					
$LP(S) = \log(0.00008377/0.0002441) = -0.4644$					
* in the case of first posit the frequency of a parti	ion in the anal cular letter in	yzed sequence the analyzed ge	we put enome		



# GLIMMER

- Gene finding program for prokaryotes (Saltzberg et. al, 1998)
- For prediction uses:
  - Start
  - Stop
  - Sequence composition
  - Interpolated Markov Models



32



#### EUKARYOTIC GENE STRUCTURE E xon 1 E xon 2 E xon 3 E xon4 DNA Intron 3 Intron2 S plice si poly A signa GGTGAG Stop codor Translatio n Branchpoin . ТА Б Л БА Л А А стаАс

#### SEARCHING FOR CODING SEQUENCES USING MARKOV CHAINS

In this case we do not want check if given sequence fragment is coding or not but we rather want to identify coding fragments in a long sequence. In most cases this is done by calculating statistics in overlapping windows.

AGTACGATATTAGCGGCAATCGTATGACTACGTCTTGCTACGTCTTCTCTCGTCTGCTCTAG



This example shows a profile for a sequence analyzed using a 120-bp window and a 10-bp step.





PROBABILITY THAT SEQUENCE IS CODING

Probability that sequence is coding is equal probability that sequence of codons is coding. Assuming independence between adjacent codons the probability that sequence is coding will be equal to the product of codon frequencies.



39

# PROBABILITY THAT SEQUENCE IS NON-CODING

If the sequence is non-coding the codon frequency will be random and each codon will be equally probable. In this case frequency for each codon will be 0.0156. This is because we have 64 codons and each of them is equally likely.

Therefore probability that the sequence is non-coding will be:

$$P(C) = F(AGG)F(AGC) = 0.0156 \times 0.0156 =$$

0.000244

40







# GENE IDENTIFICATION PROGRAMS

- ★ The first generation of programs was designed to identify approximate locations of coding regions in genomic DNA (e.g. GRAIL). These methods could not accurately predict precise exon location.
- \* The second generation (e.g. MZEF, SORFIND, and Xpound) combined splice signals and coding region identification but did not attempt to assemble predicted exons into complete genes.
- \* Third generation (GeneID, GeneParser, GenLang, FGENES) predicted entire gene structures but their performance was rather poor. One of problems was the assumption that the input sequence contains complete genes.
- ★Fourth generation of programs is represented by GENSCAN or TWINSCAN. With improved accuracy and less restricted requirements (e.g. allow partial genes) these programs are considered to be the best and are widely used in large-scale genomes analysis.

44

#### CLASSES OF GENE PREDICTION METHODS

#### Sequence similarity based

- $\cdot \ensuremath{\overset{\bullet}{\sim}}$  BLAST can be used for aligning ESTs or proteins to the genomic sequence
- PROCRUSTES and GenWise use global alignment of homologous protein to genomic sequence
- > The biggest limitation to this type of approaches:
- $\cdot \ensuremath{\succeq}$  only about half of genes being discovered have significant similarity to genes in the database
- $\cdot$  genes with very limited expression may never be discovered

#### Model based

#### Limitations of these approaches:

- Newly sequenced genomes very often lack large enough samples of known genes to estimate model parameters
- Seed to be retrained as the number of available genes is growing
  Genes of less typical structure or having rare signals may not be discovered

45

# GENSCAN

- Burge, C. and Karlin, S. (1997) Prediction of complete gene structures in human genomic DNA. J. Mol. Biol. 268, 78-94.
- Search for general and specific compositional properties of distinct functional units in eukaryotic genes
- $\overleftarrow{\begin{smallmatrix} & \end{smallmatrix}}$  General fifth-order Markov model of coding regions
- · ⊱ Analyzes both DNA strands
- Sequences may contain multiple and/or partial genes
- http://genes.mit.edu/GENSCAN

46





#### PREDICTION PROGRAMS PERFORMANCE

37 genes were tested, 16 of them (43%) were confirmed. At the exon level 159 exons were predicted and 58 (36%) were found to be real.

	predicted exons	specificity	sensitivity
MZEF	34	0.51	0.56
GRAIL	11	0.48	0.19
GENSCAN	52	0.46	0.91
FGENES	45	0.37	0.75

49



50

et al. Gene 284: 203-21;

I. Makalowska



51



52

#### GENE FINDING STRATEGIES

- $\cdot$  Search for conserved regions
- · & Presence of ORF
- •≽•Codon usage
- & Splice sites
- **⊱** Polyadenylation signal
- Similarity search
- $\cdot$  Presence of regulatory elements



#### WHY IS PROMOTER PREDICTION DIFFICULT?

- $\cdot \not \succ \cdot$  Promoter needs additional regulatory elements
- Transcription may be activated or repressed by many regulatory proteins
- ★ Transcriptional activators and repressors act very specifically both in terms of the cell type and point in the cell cycle
- $\cdot \overleftarrow{\hspace{-.05cm}\sim}$  Not all regulatory factors have been characterized





#### 57

#### EUKARYOTIC PROMOTERS

56

- \* Three types of RNA polymerase (I, II, III), each binding to various kinds of promoters
- Polymerase II transcribes genes coding for proteins
- \* Core Promoter most have TATA box that is centered around position -25 and has the consensus sequence: 5'-TATAAAA-3'
- Several promoters have a CAAT box around -90 with the consensus sequence: 5'-GGCCAATCT-3'
- · ≽ promoters for "housekeeping" genes contain multiple copies of a GCrich element that includes the sequence 5'-GGGCGG-3'
- Proximal Promoter Regions transcription factor binding regions within ~200 bp of the Core Promoter
- -> Enhancers transcription factor binding regions that can act to regulate transcription from the core promoter even from many kilobases away from the core promoter

58



#### **CISTER : CIS-ELEMENT CLUSTER FINDER**

- -> Detects cis-elements clusters by using Hidden Markov Model
- $\cdot$  For each element uses separate matrix with frequencies of each nucleotide in each position; user can input matrix for elements not included in the basic option
- User can specify:
  - .> distance between neighboring cis-elements within a cluster
  - $\cdot$  number of cis-elements in the cluster
  - · ★ distance between clusters







I I I I I I

61

### EXAMPLE OF MATRIX

Sequences of experimentally identified elements are aligned and frequencies in each position are calculated

P1 P2 P3 P4 P5 P6	NA	AML-1a					
TGTGGT	XX DE	runt-fa	ctor AM	L-1			
TCCCCT	XX						
100001	BF	T02256;	AML1a;	Species:	human,	Homo	sapiens.
TGTGGT	PO	А	с	G	т		
	01	1	0	0	4	т	
AGTGGT	02	0	0	5	0	G	
AUTOUT	03	0	1	0	4	т	
TOTOOO	04	0	0	5	0	G	
IGIGGC	05	0	0	5	0	G	
	06	0	1	0	4	т	

62

# HTTP://ZLAB.BU.EDU/ ~MFRITH/CISTER.SHTML cister : Cis-element Cluster Finder cister : Cis-element Cluster Finder Image: Cister : Cis-element Cluster Finder Cister : Cis-element Cister : Cis-element Cister : Cis-element Cister : Cister : Cis-element Cister : Cister

AND / OR upload cis-elements from a file:

63



64



# GENE EXPRESSION ANALYSIS - NCBI REPOSITORY

Gene Expression Omnibus GEO is a public functional genemics data repository supporting sequence-based data en accepted. Tools are provided to help u gene expression profiles.	MMME-compliant data submissions. Array- and sens query and download experiments and curated		Keyword or GED Accession Search
Getting Started	Tools	Browse Cont	ent
Overview	Search for Studies at GEO DataSets	Repository Brows	er .
FAQ	Search for Gene Expression at GEO Profiles	DataSets:	3848
About GEO DataSets	Search GEO Documentation	Series: 🔝	62532
About GEO Profiles	Analyze a Study with GEO2R	Platforms:	15124
About GEO2R Analysis	GEO BLAST	Samples:	1616773
How to Construct a Query	Programmatic Access		
How to Download Data	FTP Site		
Information for Submitters			
Login to Submit	Submission Guidelines	MIAME Standard	8
	Update Guidelines	Citing and Linking	to GEO
		Guidelines for Re	viewers
		GEO Publications	1



#### **USE OF RNASeq** Sequencing Task Technique approach RNASeq, small RNASeq Measuring of 36 bp single-end reads (tags counting) gene expression Measuring of RIP Seq 36 bp single-end CLIP Seq gene interactions reads (tags counting) 100 nt paired-end reads (sequence Genome RNASeq annotation determinations) 100 nt paired-end Alternative reads (sequence determinations) RNASeq splicing analysis

68













 <u>To proteinome</u> "Integrated" Transcriptome Analysis RNA Seq (Cytoplasm) AAAA Translational Control 0 tion) RNA Seq (Splicing) Post-transcriptional RNA Seq (Nuclear) egulations BS Se (DNA methy Transcriptome Sec ( AAAA mRM C DNasel Seq (Open chromatin) ATAC Seq (Open 00000 Cytoplasm 00 (TFBS f for GTF) 3C/HiC S (Higher sctru Genome Transcriptional regulations NGS as common platform

75

Bioinformatic tools for RNA Seq				
Purpose	Software	URL		
Mapping	BWA	http://bio-bwa.sourceforge.net/		
	Bowtie2	http://bowtie-bio.sourceforge.net/bowtie2/ index.shtml		
	TopHat2	http://tophat.cbcb.umd.edu/		
Expression	Cufflinks	http://cufflinks.cbcb.umd.edu/		
	Cuffdiff	Same as above		
	DEseq	http://bioconductor.org/packages/release/ bioc/html/DESeq.html		
Fusion genes	TopHat-fusion	http://tophat.cbcb.umd.edu/ fusion_index.html		
	deFuse	http://compbio.bccrc.ca/software/defuse/		
	SOAPfuse	http://soap.genomics.org.cn/soapfuse.html		
Assemble	Trans-Abyss	http://www.bcgsc.ca/platform/bioinfo/ software/trans-abyss		
	Trinity	http://trinityrnaseq.sourceforge.net/		
Viewer	UCSC Genome Browser	http://genome.ucsc.edu/cgi-bin/hgGateway		
	IGV	https://www.broadinstitute.org/igv/home		

76

#### BIOINFORMATICS CREED

- Remember about biology
- Do not trust the data
- Use comparative approach
- Use statistics
- Know the limits
- Remember about biology!!!

