


Sequence Analysis


BBSI 2006: Lecture #($\chi+1$)

Takis Benos (2006)



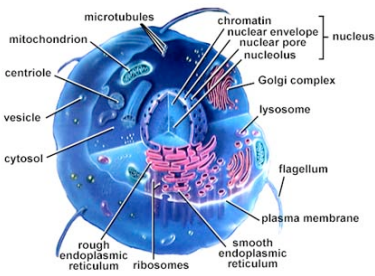

BBSI 2006 26-MAY-2006 © 2006 P. Benos

Molecular Genetics 101



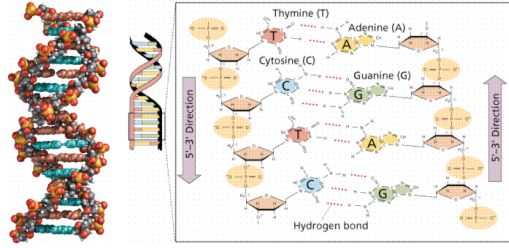
BBSI 2006 26-MAY-2006 © 2006 P. Benos

Cell's internal world

BBSI 2006 26-MAY-2006 © 2006 P. Benos

DNA - Chromosomes - Genes



BBSI 2006 26-MAY-2006

© 2006 P. Benos

What is a "gene"?

- We cannot define it (but we know it when we see it...)
- A loose definition:

"Gene" is a *DNA/RNA information unit* that is able to perform a function in a cellular environment

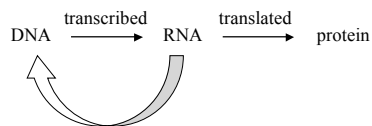


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Protein coding genes

Central Dogma:



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Open Reading Frames (ORFs)

```

aatagcgaat ttccaacga caaaagctaa atategcaaa aacctcagta aaaatcttgc 60
tggagctatt attgctaagt aacatttacc ccctgaagtt aatggatcaa tcaagagaga 120
tgtgggctgt aATGaatcgt cttattgaat taacaggttg gatcgttctt gtcgtttcag 180
tcattcttct tggcgtggcg agtcacattg acaactatca gccacctgaa cagagtgtct 240
cggtaacaaca caagTAAgct ctgcacttgt ggagcgacat gctgcccgtc cgggtgcgatg 300
ttttcaactg toggatatta aaccaggaat ttattatctt gttcgatggt gtaataaaa 358
    
```



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Open Reading Frames (ORFs)

```

aatagcgaat ttccaacga caaaagctaa atategcaaa aacctcagta aaaatcttgc 60
tggagctatt attgctaagt aacatttacc ccctgaagtt aatggatcaa tcaagagaga 120
tgtgggctgt aATGaatcgt cttattgaat taacaggttg gatcgttctt gtcgtttcag 180
M N R L
tcattcttct tggcgtggcg agtcacattg acaactatca gccacctgaa cagagtgtct 240
K stop
cggtaacaaca caagTAAgct ctgcacttgt ggagcgacat gctgcccgtc cgggtgcgatg 300
ttttcaactg toggatatta aaccaggaat ttattatctt gttcgatggt gtaataaaa 358
    
```



MNRLIELTGWIVLVVSVILLGVASHIDNYQPPEQSASVQHK



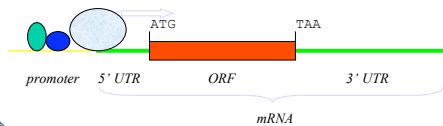
BBSI 2006 26-MAY-2006

© 2006 P. Benos

Gene's characteristics

```

aatagcgaat ttccaacga caaaagctaa atategcaaa aacctcagta aaaatcttgc 60
tggagctatt attgctaagt aacatttacc ccctgaagtt aatggatcaa tcaagagaga 120
tgtgggctgt aATGaatcgt cttattgaat taacaggttg gatcgttctt gtcgtttcag 180
tcattcttct tggcgtggcg agtcacattg acaactatca gccacctgaa cagagtgtct 240
cggtaacaaca caagTAAgct ctgcacttgt ggagcgacat gctgcccgtc cgggtgcgatg 300
ttttcaactg toggatatta aaccaggaat ttattatctt gttcgatggt gtaataaaa 358
    
```

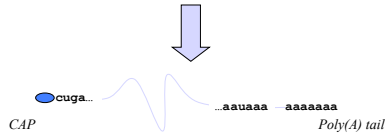


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Transcription

cugaagu aauggauca ucaagagaga 120
 ugugggcugu aAUGaaucg cuuauugaau uaacagguug gaucguucuu gucguuucag 180
 ucauuucuu uggcuggcg agucacauug acaacuauca gccaccugaa cacagugucuu 240
 cguacaaca caagUAAgcu cugcacuugu ggagcgacau gcugcccguc cgggugcaug 300
 uuuucacuug ucggauaua aaccaggaau uuauuucuu guucgauguu guauauaa 358



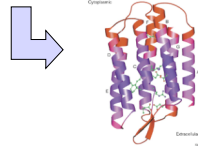
BBSI 2006 26-MAY-2006

© 2006 P. Benos

Translation

cuga... aauaaa aaaaaa

MNRLLELGGWIVLVVSVLLGVASHIDNYQPPEQSASVQHK

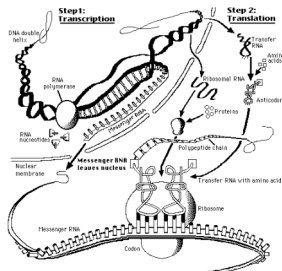


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Protein coding genes (cntd)

PROTEIN SYNTHESIS

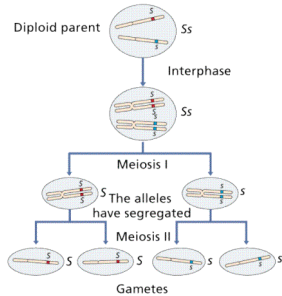


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Source: <http://www.eme.manicopa.edu/faculty/farabee/BIOBK/BioBookPR.html>
 OTS/Nr.html

Inheritance - genetic differences



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Molecular evolution

Two species will acquire mutations proportionally to their divergence time. However:

- all proteins do not change in the same pace
- a given protein does not necessarily change in the same pace throughout time
- different parts of the same protein change at different paces



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Molecular evolution (cntd)

Human (C11A_HUMAN; P05108) vs. *Drosophila* (C11A_PIG; P10612)

```

Query: 1  MLAKGLPPRSVLKVGQVTFLSAPREGLRRLVPTGEGAGISTSPRPFNFIPSPGDNGL 60
Sbjct: 1  MLAKGLALRSVLRKGCDFLSAPRECPGHPFRVGTGGCCISTATPRFSEIFSPGDNGL 60
Query: 61  NLYHFWRPCTHKVHLHVONQKQVGPVYREKLGNVSEVVIDPEDVALLFKSGPNPFR 120
Sbjct: 61  NLYHFWRKGTQRKHHVONQKQVGPVYREKLGNVSEVVIDPEDVALLFRFGPNPFR 120
Query: 121  FLIPFWAVHQYQRFVGLLKKSAAMKRDVALNOEVAPEATKNFLPLLDVSRDFVS 180
Sbjct: 121  YNIPFWAVHQYQKVFVGLLKKSGAMKRDVLRNTEVMAPEAKNFIPLLDVSRDFVS 180
Query: 181  VLHRRIKFAGSGNYSDISDLERFAFESITNVIFGEROGLMEEVNPEAQRFDAIVQM 240
Sbjct: 181  VLHRRIKQSGGKFGSDIREDLRFAPESITNVIFGERLGLMEEIVPEAQRFDAIVQM 240
Query: 241  FHTSVPMNLPPDLFLFRKTKTKRDHVAAMDVIFSKADIYTONFYWELRQKGGVHHYRG 300
Sbjct: 241  FHTSVPMNLPPDLFLFRKTKTKRDHVAAMDVIFNKAERYTONFYWDLRRKRE-FNNYFG 299
Query: 301  MLYRLGDSKMSFEDIRANVTEMLAGVDITSMTLQWHLVEMARNLKVQDMLRAEVLAR 360
Sbjct: 300  ILYRLGNDKLLSEVRANVTEMLAGVDITSMTLQWHLVEMARLNVEMLREVLINAR 359
Query: 361  HQAGDMATMLQVLLKASIKETRLRHPISVTLQRVYLVNDLVRDYMIPAKTLVQVAIV 420
Sbjct: 360  HQAGDTSMLQVLLKASIKETRLRHPISVTLQRVYLVNDLVRDYMIPAKTLVQVAIV 419
Query: 421  ALGHEPTFFPDPENFDPTRWLSKDNITYFRNLGFGWGVRCQGLGRVIAELEMFLINML 480
Sbjct: 420  AMGRDPAFFSNPQGFPTRWLGERDLHFNRNLGFGWGVRCQVGRVIAELEMFLIHIL 479
Query: 481  ENFVVEIQLSDVGTINLIMPEKISFTFWFNGEATQ 520
Sbjct: 480  ENFVVEIQLSDVGTINLIMPEKISFTFWFNGEATQ 519

```

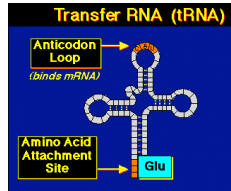


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Non-coding genes

- tRNA
- ribosomal RNA
- snoRNA
- microRNA
- etc



Source: <http://www.emc.maricopa.edu/faculty/farabee/BIOBK/BioBookPROTSYn.html>



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Elements of Probability Theory (with examples)



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Outline

- Conditional Probabilities
- Markov Chains
- Hidden Markov Models
- Information measures



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Probabilities

Definition:

$$P(x) = \frac{\text{\# favourable outcomes (x)}}{\text{total \# possible outcomes}}$$



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Conditional Probabilities

Definition:

$$P(x|A) = \frac{\text{\# favourable outcomes (x) given A}}{\text{total \# possible outcomes given A}}$$



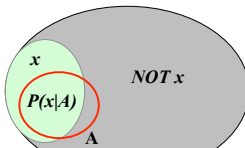
BBSI 2006 26-MAY-2006

© 2006 P. Benos

Conditional Probabilities (cntd)

Example:

- *outcome x*: tomorrow I'll go hiking
- *event A*: tomorrow will be sunny



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Conditional Probabilities (cntd)

- Joint probability: $P(X, Y) = P(X|Y) P(Y)$

- If $P(X|Y) = P(X)$ then X, Y *independent*

$$P(X, Y) = P(X) P(Y)$$

- Marginal probability:

$$P(X) = \sum_Y P(X, Y) = \sum_Y P(X|Y)P(Y)$$



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Conditional Probabilities (cntd)

- Bayes' theorem

$$P(X | Y) = \frac{P(Y | X)P(X)}{P(Y)} = \frac{P(Y | X)P(X)}{\sum_x P(Y | X)P(X)}$$

- Posterior probabilities are the compromise between data and prior information.



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Bayes: Application-1

- Problem (from Durbin *et al.*, 1998):

A rare genetic disease is discovered with population frequency one in 1 million. An extremely good genetic test is 100% sensitive (always correct if you have the disease) and 99.99% specific (false positive rate 0.01%). Will you be willing to take such a test?

- *Hint*: What is the probability that you have the disease, if the test is positive?



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Bayes: Application-1 (cntd)

- Answer:

$$\begin{aligned} P(D | +) &= P(+ | D) P(D) / P(+) = \\ &= 1.0 * 10^{-6} / [1.0 * 10^{-6} + 10^{-4} * (1 - 10^{-6})] = \\ &= 0.0099 \end{aligned}$$



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of Bayes-2

- Problem:

Given a set of transmembrane proteins with specified membrane domains of length L (*training set*), can you develop a probabilistic model that predicts which parts of a new transmembrane protein are likely to be membrane domains?



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of Bayes-2 (cntd)

- Solution:

- Suppose that we suspect that the amino acid frequencies differ between membrane and non-membrane regions.
- Using the *training set*, calculate the probabilities, $P(a_i|D)$, that each amino acid a_i is part of a membrane domain (D). Also, using the non-membrane parts, calculate the corresponding probabilities, $P(a_i|\text{not } D)$.



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of Bayes-2 (cntd)

- Solution (cntd):
 - Divide the new protein into segments.
 - Using Bayes theorem, calculate the posterior probability of each segment being a membrane domain using the $P(a_i|D)$.


$$P(X|M); M := \arg \max_M \frac{P(M|D)P(D)}{\sum P(M|d)P(d)}$$



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Markov chains

- What is a Markov chain?
- Markov chain of order n is a stochastic process of a series of outcomes, in which the probability of outcome x depends on the state of the previous n outcomes.



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Markov chains (cntd)

- Markov chain (of first order):
$$\begin{aligned} P(x) &= P(X_L, X_{L-1}, \dots, X_1) = \\ &= P(X_L | X_{L-1}, \dots, X_1) P(X_{L-1} | X_{L-2}, \dots, X_1) \dots P(X_1) = \\ &= P(X_L | X_{L-1}) P(X_{L-1} | X_{L-2}) \dots P(X_2 | X_1) P(X_1) = \\ &= P(X_1) \prod_{i=2}^L P(X_i | X_{i-1}) \end{aligned}$$
- Transition probabilities: $P(X_i | X_{i-1})$



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of Markov chains

- Problem (from Durbin *et al.*): CpG islands

Given two sets of sequences from the human genome, one with CpG islands and one without, can you calculate a model that can predict the CpG islands?



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of Markov chains (cntd)

- Solution:

	A	C	G	T		A	C	G	T
A	0.180	0.274	0.426	0.120	A	0.300	0.205	0.285	0.210
C	0.171	0.368	0.274	0.188	C	0.322	0.298	0.078	0.302
G	0.161	0.339	0.375	0.125	G	0.248	0.246	0.298	0.208
T	0.079	0.355	0.384	0.182	T	0.177	0.239	0.292	0.292

	A	C	G	T
A	-0.740	-0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679

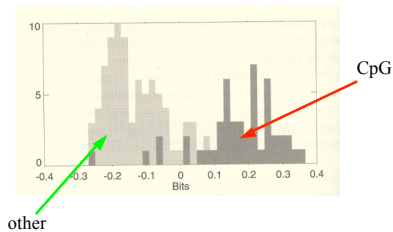


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of Markov chains (cntd)

- Histogram of scores (CpG islands):



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Hidden Markov Models

- What is a Hidden Markov Model?
- A Markov process in which the probability of an outcome depends also in a (hidden) random variable (state).
- *Transition* probability: the probability of reaching a state given the previous state.
- *Emission* probability: the probability of an outcome given the state.



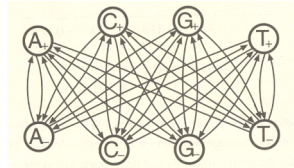
BBSI 2006 26-MAY-2006

© 2006 P. Benos

Hidden Markov Models (cntd)

- Graphical representation of the HMM:

CpG islands
(transition probabilities)



- *Question:* Where is the Markov process here?

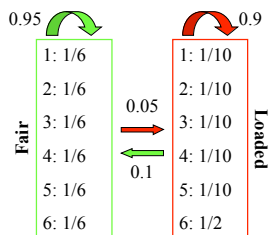


BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of HMMs

- Problem (from Durbin *et al.*): dishonest casino



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of HMMs (cntd)

- Problem (from Durbin *et al.*): dishonest casino

Given (1) the previous model and (2) a series of die rolls ($x_i, i=1, \dots, L$), can we predict which of the rolls are coming from the fair and which from the loaded die?

- Question: What is “hidden” here?



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Application of HMMs (cntd)

- Answer: YES
 - Viterbi algorithm (best path)
 - Forward-backward algorithm (probability of state k in outcome x_i)



BBSI 2006 26-MAY-2006

© 2006 P. Benos

HMMs: Viterbi algorithm

- Viterbi predictions: 300 rolls of die

```
rolls 315116246446644245311321631164152133625144543631656626566666
Die FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

rolls 65116645313265124563666463163666316232645236266666625151631
Die LLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi LLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

rolls 22255441666566563564324364131513465146353411126414626253356
Die FFFFFFFFFLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

rolls 36616366646623253441366166163252624622526525226643535336
Die LLLLLLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi LLLLLLLLLLFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF

rolls 23312162536441443233516324363366562466662632666612355245242
Die FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
Viterbi FFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFFF
```



BBSI 2006 26-MAY-2006

© 2006 P. Benos

HMMs in biology

- General comments:
 - Usually the structure of the model is unknown
 - The transition and emission probabilities are calculated based on trusted training set(s) and the postulated model
 - Predictions are based on the Viterbi or the forward-backward algorithm, depending on the question asked



BBSI 2006 26-MAY-2006

© 2006 P. Benos

Information measures

- Definitions:

- Entropy:

$$H(P) = \mathbf{E}(-\log P) = -\sum_{i=1}^n p_i \log p_i$$

- Relative Entropy:

$$H(P, Q) = \sum_{i=1}^n p_i \log \frac{p_i}{q_i}$$

- Mutual Information:

$$M(X, Y) = \sum_{i,j} P(x_i, y_j) \log \frac{P(x_i, y_j)}{P(x_i)P(y_j)}$$



BBSI 2006 26-MAY-2006

© 2006 P. Benos
