

PONENCIA # 4

TOOLBOX DE BIOINFORMATICA: ENTORNO DE SOFTWARE INTEGRADO PARA EL GENOMA Y ANÁLISIS PROTEÓMICO

JAIRO PERTUZ CAMPO

Físico, Instructor y asesor de MATLAB, E.MAILS: jpertuz@udem.edu.co, pertuzjairo@yahoo.es

SOLON PINO G.

D. SOFTWARE PARA INGENIERIA

COMPONENTES ELECTRONICAS LTDA

1. PRESENTACIÓN DEL BIOINFORMATICS TOOLBOX

Componentes Electrónicas Ltda Distribuidor Autorizado en Colombia de **The MathWorks Inc.** anunció el 6 de Noviembre la disponibilidad de su **nuevo Bioinformatics Toolbox** para MATLAB ®. El producto facilita el análisis de datos complejos de bioinformática, de las sucesiones de análisis genómico de datos en micro-arreglos, y velocidades de desarrollo y despliegue de algoritmos. Como el Toolbox se construye en MATLAB ®, los especialistas en bioinformática pueden confiar en un análisis de datos probado y un ambiente de desarrollo de algoritmos que les permiten responder eficiente y rápidamente, con menos tareas de codificación.

"MATLAB ® proporciona una abundancia de herramientas que son críticas al planear el desarrollo de un algoritmo. El Bioinformatics Toolbox se suma naturalmente a MATLAB ® manteniendo una interfaz sin limitaciones y accediendo con él a los almacenes de datos del genoma esenciales y de dominio público, y proporcionando los constructores de bloques necesarios para procesar los datos "genómicos" dijo Bahram Ghaffarzadeh Kermani, Ph. D., del personal científico de Illumina, Inc. en San Diego.

El Bioinformatics Toolbox se diseña para ayudar a los usuarios a entender y visualizar cantidades grandes de datos de la investigación asociada con las aplicaciones de la tecnología biológica o biotecnología tradicionales, en un fragmento del tiempo. Con este último Toolbox The MathWorks está entregando el poder y versatilidad de su ambiente informático técnico integrado directamente a la bio-tecnología y las industrias farmacéuticas. Aprovechando la funcionalidad y riqueza de la programación en MATLAB ®, el Bioinformatics Toolbox proporciona un ambiente abierto, extensible a científicos y a bioinformáticos para el desarrollo de algoritmos, manipulación de este tipo de datos y habilita la comunicación fácil de aplicaciones por las diferentes secciones de la investigación. Como resultado, los bioinformáticos pueden usar el toolbox para enfocar los esfuerzos en el centro de su trabajo - la investigación y análisis - sin los riesgos asociados con usar programas o software dispares.

:

Entre sus numerosos rasgos y capacidades, el Bioinformatics Toolbox proporciona el acceso a archivos del genoma en formatos normales, los bancos de datos basados en la web como GenBank y PIR, y las fuentes de los datos en línea. El toolbox también ofrece las rutinas especializadas para visualizar los datos de microarray (micro-arreglos), incluso las cajas de gráfico, los gráficos I-R y los mapas espaciales de calor.

Los rasgos importantes adicionales del Bioinformatics Toolbox incluyen:

- Archivos y expresiones en formato comprensibles en genética, genómica y proteómica.
- Conversión, adaptación y estadísticas de secuencias del genoma y proteínas.
- Funciones de análisis proteínico.
- Gráficas de puntos, gráficas de grupos, gráficas de sectores y otras representaciones gráficas de datos genómicos y proteómicos.

.

" Los bioinformáticos han tenido que invertir mucho tiempo en matemática de programación y algoritmos de estadística en un horario corto, tradicionalmente," dijo Kristen Amuzzini, gerente del área de biotecnología, farmacéutica e industria médica de, The MathWorks. Y continua : "La combinación de los productos de MATLAB ® de hoy y el nuevo Bioinformatics Toolbox entrega las herramientas que ellos necesitan para analizar los datos grandes y como resultado da elementos que identifican los puntos donde potencialmente se necesita desarrollar un remedio rápida y eficazmente."

No dude en preguntarnos por el nuevo **Toolbox de Bioinformática**, es parte de una solución agil y efectiva que puede representarle mucho ahorro de tiempo y dinero, en biotecnología, farmacéutica y medicina.

Si lo desea Usted puede abrir en www.mathworks.com el seminario

2. EJEMPLO DE ANALISIS DE SECUENCIAS DEL GENOMA HUMANO:

SEQSTATSDemo Example of sequence statistics with MATLAB

This demonstration looks at some statistics about the DNA content of the human mitochondrial genome.

Contents

- Introduction
- Composition of the mitochondrial genome.
- Exploring the Open Reading Frames (ORFs)
- Extracting and analyzing the ND2 protein

Introduction

Mitochondria are generally the major energy production center in eukaryotes. The Genome repository at the NCBI contains more interesting information about the human mitochondrial genome.

```
web(http://www.ncbi.nlm.nih.gov/)
```

The consensus sequence of the human mitochondria genome has accession number NC_001807. The whole GenBank entry is quite large and this example only uses the nucleotide sequence, so you can use the **getgenbank** function with the 'SequenceOnly' flag to read just the sequence information into the MATLAB workspace.

```
mitochondria = getgenbank('NC_001807', 'SequenceOnly', true);
```

If you don't have a live web connection, you can load the data from a MAT-file using the command

```
% load mitochondria % <== Uncomment this if no internet connection
```

The MATLAB **whos** command gives information about the size of the sequence.

```
whos mitochondria
Name          Size            Bytes   Class
mitochondria    1x16571        33142  char array

Grand total is 16571 elements using 33142 bytes
```

You will use some of the sequence statistics function in the Bioinformatics Toolbox to look at various properties of this sequence. You can look at the composition of the nucleotides with the **ntdensity** function.

```
ntdensity(mitochondria)
```

PONENCIA # 4

Composition of the mitochondrial genome.

This shows that the genome is A-T rich. You can get more specific information with the **basecount** function.

```
basecount(mitochondria)
ans =
A: 5113
C: 5192
G: 2180
T: 4086
```

These are on the 5'-3' strand. You can look at the reverse complement using the **seqrcomplement** function.

```
basecount(seqrcomplement(mitochondria))
ans =
A: 4086
C: 2180
G: 5192
T: 5113
```

As expected, the base counts on the reverse complement strand are complementary to the counts on the 5'-3' strand.

You can use the chart option to **basecount** to display a pie chart of the distribution of the bases.

```
figure
basecount(mitochondria, 'chart', 'pie');
```

Now look at the dimers in the sequence and display the information in a bar chart using **dimercount**.

```
figure
dimercount(mitochondria, 'chart', 'bar')
ans =
AA: 1594
AC: 1495
AG: 801
AT: 1223
CA: 1536
CC: 1779
CG: 439
CT: 1438
GA: 615
GC: 716
GG: 427
GT: 421
TA: 1368
TC: 1202
TG: 512
TT: 1004
```

PONENCIA # 4

You can look at codons using **codoncount**. The function **dimercount** simply counts all adjacent nucleotides; however **codoncount** counts the codons on a particular reading frame. With no options, the function shows the codon counts on the first reading frame.

```
codoncount(mitochondria)
AAA - 172    AAC - 157    AAG - 67    AAT - 123
ACA - 153    ACC - 163    ACG - 42    ACT - 130
AGA - 58     AGC - 90     AGG - 50     AGT - 43
ATA - 132    ATC - 103    ATG - 57     ATT - 96
CAA - 166    CAC - 167    CAG - 68     CAT - 135
CCA - 146    CCC - 215    CCG - 50     CCT - 182
CGA - 33     CGC - 60     CGG - 18     CGT - 20
CTA - 187    CTC - 126    CTG - 52     CTT - 98
GAA - 68     GAC - 62     GAG - 47     GAT - 39
GCA - 67     GCC - 87     GCG - 23     GCT - 61
GGA - 53     GGC - 61     GGG - 23     GGT - 25
GTA - 61     GTC - 49     GTG - 26     GTT - 36
TAA - 136    TAC - 127    TAG - 82     TAT - 107
TCA - 143    TCC - 126    TCG - 37     TCT - 103
TGA - 64     TGC - 35     TGG - 27     TGT - 25
TTA - 115    TTC - 113    TTG - 37     TTT - 99
```

Using a loop you can also look at all the other reading frames:

```
for frame = 1:3
figure
subplot(2,1,1); codoncount(mitochondria,'frame',frame,'figure',true);
title(sprintf('Codons for frame %d',frame));
subplot(2,1,2);
codoncount(mitochondria,'reverse',true,'frame',frame,'figure',true);
title(sprintf('Codons for reverse frame %d',frame));
end
```

Exploring the Open Reading Frames (ORFs)

In a nucleotide sequence an obvious thing to look for is if there are any open reading frames. The function **seqshoworfs** can be used to visualize ORFs in a sequence. Note: In the HTML tutorial only the first 7500 bases of the first reading frame are shown, however when running the demo you will be able to inspect the complete mitochondrial genome with the aid of the **Help Browser**.

```
seqshoworfs(mitochondria);

Frame 1

000001      gatcacaggcttatcacccattaaaccactcacggagacttccatgcattttgtatcg
000065      tgggggggtgtgcacgcgatagcattgcgagacgtggagccggagcaccctatgtcgactatc
000129      tgtcttgattcctgcctatttattatcgcacctacgttcaatattacaggcgaacat

http://matlab.udes.edu.co
http://es.geocities.com/matlab\_colombia/diamatlabnov3.html
http://www.compelect.com.co/FormularioDiaMATLAB.html
```

PONENCIA # 4

000193 acctactaaaagtgttaattaattaatgctttaggacataataacaattgaatgtctgc
 000257 acagccgcttccacacagacatcataacaaaaatttcaccaaaccggccctcccccgct
 000321 tctggccacagcacctaaccacatctctgcaccaaccggaaaacaaagaaccctaaccaggcc
 000385 taaccagattcaaaattttatctttaggcgtatgcactttacagtcaccccaactaaca
 000449 cattattttcccccactccctactaataatctcatcaataacaaccggccatccacc
 000513 cagcacacacaccgcttaaccctacccgaaccaaccaaaacccaaagacacccca
 000577 cagtttatgttagcttaactcctcaagcaataactgaaaatgtttagacgggctcacatcacc
 000641 ccataaacaataggttgcctagcctttattagcttttagtaagattacacatgcaagc
 000705 atccccgtttcagtgagttcacccctctaattaccacgatcaaaaggacaagcatcaagc
 000769 cagcaatgcagctcaaaacgcttagcctagccacacccacgggaaacagcagtgattaaact
 000833 ttagcaataaaacgaaagttactaagctataactaaccgggttggtaattcgtgccagc
 000897 caccgcgtcacacgatataagaagccgtaaagagtgttttagatcaccc
 000961 cctccccaataaagctaaaactcacctgagttgtaaaaaactccagttgacacaaaatagacta
 001025 cggaaagtggcttaacatatctgaacacacaatagctaagacccaaactgggattagatacc
 001089 actatgcttagccctaaacctcaacagttaaatcaacaaaactgctccagaacactacgagc
 001153 cacagcttaaaactcaaaggacctggcggtcatatccctctagaggagccgttctgtaa
 001217 tgatataaccccgatcaacccctaccaccttgcagcttatataccgcatctcagcaaac
 001281 cctgtatgaaggctacaagaatgcaagtgtaaccacgtaagacgcttaggtcaaggtgttagccc
atgaggtggcaagaaatggctacatccatcccaactacgatagcccttatgaaact
 001345 taagggtcgaagggtggatttagcagtaactgagagtagagtgttagtgaacaggggccctga
 001409 agcgcgtacacaccgcgtcaccctctcaagtataacttcaaggacatttactaaaacccc
 001473 tacgcatttatatacggagacaagtcgttaacatggtaagtgtactggaaagtgcacttggacg
 001537 aaccagagtgttagcttaacacaaaagcacccactacacttaggagattcaacttaacttgac
 001601 cgctctgagctaaacctagcccaaccactccactaccagacaacccctagccaaacc
 001665 atttacccaaataaagttataggcgatagaaattgaaaccctggcgcaatagatatacgccaa
 001729 gggaaagatgaaaaattataaccaagcataatatacgaaaggactaaccctatacctctgc
 001793 atatgaattactagaaataactttgcaaggagagccaaagctaaagaccccgaaaccagacg
 001857 ctacctaagaacagctaaagagcacccgtctatgttagctaaatagtgaaagatttataagg
 001921 tagaggcgacaaacctaccgcgctgggtatagctgggttcaagatagaatcttagtcaac
 001985 tttaatttgcacagaaccctctaaatccctgttaatttactgttagtccaaagagggaa
 002049 cagctcttggacacttaggaaaaacccctgttagagagtagaaaaatttaacccatagtagg
 002113 cctaaaagcagccaccaattaagaagcgttcaagctcaacaccactacccctaaatccca
 002177 acatataactgaactcctcacaccaattggaccaatctatcaccctatagaagaactaatgt
 002241 agtataagtaacatgaaaacattctccctcgacataaccgcgtcagatcaaaacactgaactg
 002305 acaattaacagcccaatatctacaatcaaccaacaagtcttaccctactgtcaacccaa
 002369 cacaggcatgtcataaggaaaggttaaaaaaagtaaaaggactcggccaaacccttacccg
 002433 tggaccatccatcacccattcttagcatcaccgtttagaggcaccgcctgccc
 002497 **tgtttaacggccgcgtaccctaaccgtgcaaaagg**tagcataatcacttgttccctaaatagg
 002561 acctgtatga~~atggctccacgagggttcagctgtcttactttaccagtgaaattgacctg~~
cccgtaagaggcggcatgacacagcaagacgagaagaccctatggagcttaatttataat
 002689 gcaaacagtacctaacaaccccacaggtctaaactaccaaaacctgcattaaaatttgc
 002753 gggcgacccctggagcagaacccaaacctccgagcgtacatgctaaagacttaccagt
 002817 caaagcg
 002881 aactactataactcaattgtaccaataacttgaccaacggaaacttaccctaggataac
 002945 gcaatccattctaggtccatatacaataaggtttacgacccgtatgttggatcaggadat
 003009 **cccgatggcagccgtattaaaggttcgttcaacgat**taagtccctacgttatctg
 003073 ttacgaccggagtaatccaggcgtttctatctacttcaaattccctccctgtacgaaagg
 003137 agagaataaggcctacttcacaaaggcccttccccc
gttaatgatatcatctcaacttagtat
tatacccacaccaccaagaacagggttgttaagatggcagagcccgtaatcgcaaaaaac
 003201 ttaaaacttacagtcaagggttcaattctttaacaacataccatggccaaaccctcta
 003265 ctccctattgttaccatctaatcgcaatggcattcttaatgcttaccgaacggaaaattct
 003329 tag
 003393 gctatataactacgcaaaggcccaacgttgcgttaggccttacgggctactacaacc
 003457 ctgc
 003521 tgacgccataaaaactctcaccacaaagagccctaaaaccgcacatctaccatacc
 003585 ctacccgcggccacccatcgcttactatgatccaccatcgatcttact
atgaaccccccctccccataccca
acccctggtcaacctcaaccttaggcctcttatttacttagccacotctagcctagccgtt
 003649 ctaatccctgtatcagggtgagcatcaaactacgcccgtatcggcgcactgc
 003713 gtggctctttaaccttccacccttaccaacacaagaacacac
 003777 ctgtattacttccacactacgagaccaaccccttgc
atgaccctggccataatatgattatctccacactacgagaccaaccccttcgac

PONENCIA # 4

003905 cttggccgaaggggagtcgcactagtctcaggctcaacatcgaaatacgcgcaggcccccttcg
003969 ccctattcttcatagccgaaatacacaacattattataataaaacaccctcaccactacaatctt
004033 ccttaggaacaacatatgacgcactctccctgaactctacacaacatatttgtcaccaagacc
004097 ctacttctaaccctccgttctt atgaattcgaacagcataccccgattccgctacgaccacaa
004161 tcatacacctctatgaaaaacttcctaccactcaccctagcattacttatatgatatgtctc
004225 cataccattacaatctccagcattcccccctaagaaaatgtctgataaaagagtt
004289 ctttgatagagtaaataataggagcttaaaccccttattttaggactatgagaatcgAACCC
004353 atccctgagaatccaaaattctccgtgccacctatcacacccatctcaaagtaaggctcgacta
004417 aataagctatcgccccatacccgaaaaatgtgttataccctccgtactaattaatcccc
004481 tggcccaacccgtcatctacttaccatcttgcaggcacactcatcacagcgctaagctcgca
004545 ctgatttttacctgagtaggcctagaaaataacatgctagttttattccagttctaaacccaa
004609 aaaataaacccctcggtccacagaagctgcacatcaagtatttctcagcgaagcaaccgcattcca
004673 taatccttctaatagctatccttcaacaataactctccggacaatgaaaccataaccaataac
004737 taccatcaataactcatcattaataatcata atgctatagcaataaaacttaggaatagcccccc
004801 tttcacttc tgagtcccagaggttaccaaggcaccctctgacatccggcctgcttcttctca
004865 catgacaaaactagccccatctcaatcatataccaaatctccctactaaacgtaagcct
004929 tctcctactctcaatcttattccatcatagcaggcagttgagggtggattaaacccaaacccag
004993 ctacgcaaaatcttagcatactcccaattacccacataggatgaataatagcagtctaccgt
005057 acaaccctaaacataaccatcttaattaacttattatattatcctaactactaccgcattct
005121 actactcaacttaaactccagcaccacgaccctactactatctgcacctgaaacaagctaaca
005185 tgactaacaccctaattccatccaccctctcccttaggaggcctgccccgctaaccggct
005249 tttgccccaaatggggcattatcgaagaattcacaaaaaacaatagccctcatcatccccaccat
005313 catagccaccatcacccctcttaacccttacttctacctacgcctaattctactccacccatcaatc
005377 acactactcccccataatctaacaacgtaaaaataatgacagttgaacataaaaaacccacc
005441 cattcctccccacactcatcgcccttaccacgcactcttctccctttataactaat
005505 aatcttataaaaaatttaggttaaatacagaccaagagccttcaaagccctcagtaagtgcatt
005569 acttaatttctgcaacagcataaggactgcaaaaacccactctgcacactgaacgcacatc
005633 ccactttaattaagctaaggccctactagaccaatggacttaaaccacaaacacttagtt
005697 cagctaagcaccctaactcaactggctcaatctacttctccgcggaaaaaaggcggga
005761 gaagccccggcagggttgaagctgcttctcgaattcaatataaaaaatcaccccgaa
005825 gctggtaaaaagaggcctaaccctgtcttagatttacagttcaatgcttactcagccatt
005889 tacctcacccccactgtatgttgcgcgaccgttactattcttacaaaaccacaaagacattt
005953 acactataccattattcggcgcatacgactggagtcctaggcacagctctaagcccttattt
006017 gagccgagctggccagccaggcaaccccttaggttaacgcaccacatctacaacgcttac
006081 agcccatgcattttaataatcttcttcatagtaataccatcataatccggaggcttggcaac
006145 tgactagttcccttaataatcggtgcggccgat atggcgccccgcataaaacaacataagct
006209 tctgactcttacccctctccactcctgctcgatctgtctatagtggaggccggagcagg
006273 aacaggttgaacagctaccctcccttagcaggaaactactcccaccctggagccctccgtagac
006337 ctaaccatcttccttacacccctagcagggtgtctcttatcttagggccatcaatttcatca
006401 caacaattatcaatataaaacccctgccataacccaaatccaaacgccccttctcgctgtatc
006465 cgtcttaatcacagcagtctacttctctatcttccctagtcgtctgcatcactata
006529 ctactaacagaccgcaaccccttacccatcttcgaccccccggagaggagacccattc
006593 tataccaacacccattctgattttccgttaccctgaagttatattcttaccaggctt
006657 cggaaataatctcccatattgtacttactactccggaaaaaaaagaaccattggatatacataggt
006721 atggctgagct atgatataattgggttccctagggtttatctgttgagcacaccatataattt
006785 cagtaggaatagacgttagacacacgagcatattcacctccgtaccataatcatcgctatccc
006849 caccggcgtcaaagtatttagctgactcgccacactccacccgaaagcaatatgatctgt
006913 gcagtgcgtcgacccttaggattcatcttcttcccgtaggtggctgactggattgtat
006977 tagcaaactcatcactagacacatcgtaactacacgcacacgtactacgtttagctcacttccacta
007041 tgtcctatcaataggagctgtatttgcattcataggaggcttcatctgatttcccttattc
007105 tcaggctacaccctagacccatccacccatattctactatcatattcatcgccgtaa
007169 atctaacttctcccacaacactttctcgccctatccggaa atgccccgacgttactcgacta
007233 cccccatgcatacaccacatgaaacatcttcatctgttaggcttattcttcttaacagca
007297 gtaatattaataatttcatgatttggagaagccctcggtcaagcgaaaaagtcttaatagtag
007361 aagaaccctccataaacctggagtgactatatggatgccccccaccctaccacacattcgaaga
007425 acccgatatacataaaaatctagacaaaaaaggaaggaatcgaaaccccccacccatcgaa
007489 caaccccatqqc

PONENCIA # 4

If you compare this output to the genes shown on the NCBI page there seem to be slightly fewer ORFs, and hence fewer genes, than expected. Vertebrate mitochondria do not use the Standard genetic code so some codons have different meaning in mitochondrial genomes. For more information about using different genetic codes in MATLAB see the help for the function **geneticcode**.

```
help geneticcode
```

GENETICCODE returns a structure containing mappings for the genetic code.

MAP = GENETICCODE returns a structure containing mapping for the Standard genetic code.

GENETICCODE(ID) returns a structure of the mapping for alternate genetic codes, where ID is either the transl_table ID from the NCBI Genetics web page (<http://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi?mode=c>) or one of the following supported names. NAME can be truncated to the first two characters of the name.

ID Name

1	<i>Standard</i>
2	<i>Vertebrate Mitochondrial</i>
3	<i>Yeast Mitochondrial</i>
4	<i>Mold, Protozoan, and Coelenterate Mitochondrial and Mycoplasma/Spiroplasma</i>
5	<i>Invertebrate Mitochondrial</i>
6	<i>Ciliate, Dasycladacean and Hexamita Nuclear</i>
9	<i>Echinoderm Mitochondrial</i>
10	<i>Euplotid Nuclear</i>
11	<i>Bacterial and Plant Plastid</i>
12	<i>Alternative Yeast Nuclear</i>
13	<i>Ascidian Mitochondrial</i>
14	<i>Flatworm Mitochondrial</i>
15	<i>Blepharisma Nuclear</i>
16	<i>Chlorophycean Mitochondrial</i>
21	<i>Trematode Mitochondrial</i>
22	<i>Scenedesmus Obliquus Mitochondrial</i>
23	<i>Thraustochytrium Mitochondrial</i>

Examples:

```
moldcode = geneticcode(4);
wormcode = geneticcode('Flatworm Mitochondrial');
```

See also NT2AA, REVGENETICCODE.

The 'GeneticCode' option to the **seqshoworfs** function allows you to look at the ORFs again but this time with the vertebrate mitochondrial genetic code. Notice that there are now two much larger ORFs on the first reading frame: One starting at position 4471 and the other starting at 5905. These correspond to the ND2 (NADH dehydrogenase subunit 2) and COX1 (cytochrome c oxidase subunit I) genes.

```
orfs = seqshoworfs(mitochondria,'GeneticCode','Vertebrate Mitochondrial',...
    'alternativestart',true)
orfs =
```

PONENCIA # 4

Start: [1x26 double]
Stop: [1x26 double]

Frame 1

```

000001 gatcacaggctatcaccctattaaccactcacgggagctctccatgcatttggtatttcg
000065 tgggggggtgtgcacgcgatagcattgcgagacgctggagccggagcaccctatgtcgca
000129 gatc
000193 acctactaaagtgtttaattaattaatgtctttaggacataataaacaattgaatgtctgc
000257 acagccgccttcacacagacatcataacaaaaatttccaccaaacc
000321 tctggccacagcacttaaacacatctgc
000385 taaccagattcaaaattttatctttaggcggatgcactttaacagt
000449 cattattttccctcccactcccataactactaatctcatcaata
000513 cagcacacacacaccgctgctaacc
000577 cacc
000641 ccataaaacaaataggttggcttagc
000705 atccccgttccagtgagttcac
000769 cacc
000833 ttagcaataaaacgaaagttaacta
000897 aactaacc
000961 ccc
001025 ctt
001089 actatgcttagccctaaac
001153 ct
001217 tgat
001281 ctgat
001345 atgggt
001409 taagggt
001473 agcgcgt
001537 tacgc
001601 acc
001665 cgt
001729 attac
001793 gggaa
001857 aat
001921 ac
001985 tagaggc
002049 tt
002113 cagct
002177 ctt
002241 acatata
002305 agt
002369 aca
002433 cacagg
002497 ttttac
002561 tgg
002625 ac
002689 cccgt
002753 gca
002817 gggc
002881 aactactata
002945 gca
003009 cccgat
003073 ttc
003137 agaga

```

PONENCIA # 4

003201 **tatacccacaccacccaagaaacagggttgttaagatggcagagcccgtaatcgataaaaac**
 003265 taaaaacttacagtcaaggttcaattctttcttaacaacataaccatggccaaccttccta
 003329 ctctcattgtaccattctaatcgcaatggcattcctaattgttaccgaacgaaaaattctag
 003393 **gtatatacaactacgcaaaggcccaacgttgtaggcccctacgggtactacaacccttgc**
 003457 tgacgccataaaactcttccaccaaagagccctaaaaccgcac **atctaccatacccttac**
 003521 **atcaccggccgaccttagtctcaccatcgcttctactatgaacccccctccccataccca**
 003585 **acccctggtcaacctcaaccttaggcctctatttattctagccaccccttagcctagccgtt**
 003649 ctcacatcctgtatcagggtgagcatcaaactaaactacgcccgtatggcgactgcgagca
 003713 **gtagcccaaacaatctcatatgaagtccatccatgttactatcaacattactaataaa**
 003777 **gtggctccttaaccttcacccttatcacaaacacaagaacaccctgttactatcgtccatc**
 003841 **atgacccttgccataatgattatctccacactagcagagaccaaccgaaccccttcgac**
 003905 ctggccgaaggggagtccgaacttagtctcagggttcaacatcgaaataccgcgcaggcccccttcg
 003969 ccct **attcttcatgccatcacaaacattataataaaaccccttaccactataatctt**
 004033 **ccttaggaacaacatatgacgcactctccctgaaactctacacaatatttgcaccaagacc**
 004097 **ctacttctaaccctccctgttcttatgtatcgacacagcatccccgattccgtacgaccaac**
 004161 **tcatacacccatctatgaaaaactcttaccactcacccttagcattactatatgatatgtctc**
 004225 **catacccatacaatctccagcatccccctcaaaacctaagaaatatgtctgataaaagagta**
 004289 ctttgatagagtaataataggagcttaaccccttatttctaggactatgagaatcgacc
 004353 atccctgaga **atccaaaattctccgtccacccatcaccccatcttcaagtaaggctagcta**
 004417 aataagctatggggccatacccccggaaaatgttggttatccctccctgtactattaatcccc
 004481 **tggcccaacccgtcatctactcttgcaggcacactcatcacagcgtaagctcgca**
 004545 **ctgatTTTtaccttgagtaggccttagaaataaacatgtcttttattccgttctaaacaaa**
 004609 **aaaataaaacccctgttccacagaagctgcatcaagtattccctcacgcaagcaaccgcacca**
 004673 **taatccttctaatacgatcttcaacaatatactctccggacaatgaaccataaccatac**
 004737 **taccaatcaatactcatcattaataatcataatggctatagcaataaaacttaggaatagcccc**
 004801 **tttcaacttcttgagttccagggtaaccaaggcaccctctgacatccggctgttcttctca**
 004865 **catgacaaaaactagccccatctcaatcatataccaaatctctccctcaactaacgtaaagcct**
 004929 **tctcctcaactctcaatcttacccatcatagcaggcaggtaggtgattaaaccaaaaaac**
 004993 **ctacgaaaaatcttagcataactcccaattaccacataggatgaataatagcagttctaccgt**
 005057 **acaaccctaaacataaccattcttaatttaactatttattatcttaactactaccgcattct**
 005121 **actactcaactaaactccagcaccacgaccctactactatctcgacccgtaaacaagctaaca**
 005185 **tgactaacacccttaatccatccaccctctcccttagggaggcctgccccctgtaccgg**
 005249 **tttgcaccaatggccattatcgaaagaattcacaaaaacaatagccatcatccccccat**
 005313 **catagccaccatcacccctttaacctctacttctacccctacgtccataatcttactccatca**
 005377 **acactactcccatatctaacaacgtaaaaataaaatgacagttgaacatacaacccaccc**
 005441 **cattccctcccacactcatgccccttaccacgctactccctacctatctccctttataacta**
 005505 **aatcttaatgaaattttaggttaataacagaccaagagccctcaagccctcagtaagttgcaat**
 005569 **acttaatttctgcaacagactaaggactgaaaacccactctgcataactgaacgcaatcag**
 005633 **ccactttaattaagctaagcccttactagaccaatggacttaaaccacacaacttagttaa**
 005697 **cagctaagcaccctaatcaactggctcaatctacttctccgcgcggggaaaaaggcg**
 005761 **gaagccccggcagggttgaagctgttctcgaaatttgcattcaatgtaaaatcacctcgga**
 005825 **gctggtaaaaaagaggcctaaccctgtcttagattacgtccaaatgttctcactcagccatt**
 005889 **tacctcaccccaactgatgttgcggaccgttacttctacaaaccacaaagacatttgg**
 005953 **acactataccattatttccgcgtatggactgttccatgttccatgttccatgttccat**
 006017 **gagccgagctggccagccaggcaacccatcttagttaacgaccacatctacaacgttacgttac**
 006081 **agcccatgcatttgcataatcttcttacccatcataatccggaggcttggcaac**
 006145 **tgacttagttcccttaataatcggtccccgatatggcgtttcccgccatcaaacaacataag**
 006209 **tctgactcttacccctctctactcctgtctcgcatgtctatagtgaggccggaggcagg**
 006273 **aacaggttgaacagtctaccctcccttagcaggactactccaccctggagctccgttagac**
 006337 **ctaacccatcttctcccttacccatcgagggtctctctatcttagggccatcaatttcatca**
 006401 **caacaattatcaataaaaaacccctgtccataacccataccacgcccccttctgtgtatc**
 006465 **cgtccctaatcacagcaggctacttctctatcttccctgttccatgtctgttccatgttccat**
 006529 **ctactaacagaccgcaacccatcaacaccaccccttccgtaccccgccggaggaggacccat**
 006593 **tataccaaacacctattctgatTTTcggttacccatgttccatgttccatgttccatgttccat**
 006657 **cggataatctccatattgttaactacttccggaaaaaaaagaaccatgttccatgttccat**
 006721 **atggtctgagctatgtatcaatttgctttaggtttatctgtgtgagcacaccatatttgc**
 006785 **cagtaggaatagacgttagacacacgagcatatttccctccgttccatgttccatgttccat**
 006849 **caccggcgtaaaagtatttagtgcactcgccacactccacccggaaagcaatatgtatctgt**

PONENCIA # 4

```

006913  gcagtgcgtcgagccctaggattcatcttctttcacggtaggtggcctactggcattgtat
006977  tagcaaactcatcaactagacatcgtaactacacgacacgtactacgtttagctcaactccacta
007041  tggcttatcaataggagctgtattgccatcataggaggcttcatctactatgttcccatttc
007105  tcaggctacaccctagacaaaacccatcgccaaatccatctactatcatattcatcggtaa
007169  atctaactttcttccacaacacttctcggtatccgaatgccccgacgttactcggaacta
007233  ccccgatgcatacaccatgaaacatcctatcatctgtaggcttattctctaacagca
007297  gtaatattaataatttcatgattttagaaagccttcgttgcgaagcgaaaagtcttaatagtag
007361  aagaaccctccataaaacctggagtgactatatggatgccccccaccctaccacacattcgaaga
007425  acccgatatacataaaaatctagacaaaaaaggaaggatcgaacccccaagctggttcaagc
007489  caacccatggc

```

Extracting and analyzing the ND2 protein

The ORF of interest starts at position 4471, the following commands can be used to find the corresponding stop codon:

```

ND2Start = 4471;
startIndex = find(orfs(1).Start == ND2Start)
ND2Stop = orfs(1).Stop(startIndex)
startIndex =

```

24

ND2Stop =

5512

Once the positions are known, MATLAB indexing can be used to extract the region of interest.

```
ND2Seq = mitochondria(ND2Start:ND2Stop);
```

If you look at the **codoncount** for this gene we see a lot of CTA and ATC codons.

codoncount(ND2Seq)			
AAA - 10	AAC - 14	AAG - 2	AAT - 6
ACA - 11	ACC - 24	ACG - 3	ACT - 5
AGA - 0	AGC - 4	AGG - 0	AGT - 1
ATA - 22	ATC - 24	ATG - 2	ATT - 8
CAA - 8	CAC - 3	CAG - 2	CAT - 1
CCA - 4	CCC - 12	CCG - 2	CCT - 5
CGA - 0	CGC - 3	CGG - 0	CGT - 1
CTA - 26	CTC - 18	CTG - 4	CTT - 7
GAA - 5	GAC - 0	GAG - 1	GAT - 0
GCA - 8	GCC - 7	GCG - 1	GCT - 4
GGA - 5	GGC - 7	GGG - 0	GGT - 1
GTA - 3	GTC - 2	GTG - 0	GTT - 3
TAA - 0	TAC - 8	TAG - 0	TAT - 2
TCA - 7	TCC - 11	TCG - 1	TCT - 4
TGA - 10	TGC - 0	TGG - 1	TGT - 0
TTA - 8	TTC - 7	TTG - 1	TTT - 8

PONENCIA # 4

For those of you who have not memorized the genetic code you can easily check what amino acids these codons get translated into using the **nt2aa** and **aminolookup** functions.

```
aminolookup('letter',nt2aa('CTA'))
aminolookup('letter',nt2aa('ATC'))
```

ans =

Leu leucine

ans =

Ile isoleucine

The **nt2aa** function converts the nucleotide sequence to the corresponding amino acid sequence. Again the 'GeneticCode' option must be used to specify the vertebrate mitochondrial genetic code.

```
ND2 = nt2aa(ND2Seq,'GeneticCode','Vertebrate Mitochondria');
```

You can get a more complete picture of the amino acid content with **aaccount**.

```
figure
aaccount(ND2,'chart','bar')
ans =
```

```
A: 20
R: 4
N: 20
D: 0
C: 0
Q: 10
E: 6
G: 13
H: 4
I: 31
L: 64
K: 12
M: 25
F: 15
P: 23
S: 28
T: 43
W: 11
Y: 10
V: 8
```

Notice the high leucine, threonine and isoleucine content and also the lack of cysteine or aspartic acid.

You can use the **atomiccomp** and **molweight** functions to find out more about the ND2 protein.

<http://matlab.udes.edu.co>
http://es.geocities.com/matlab_colombia/diamatlabnov3.html
<http://www.compelect.com.co/FormularioDiaMATLAB.html>

PONENCIA # 4

```
atomiccomp(ND2)
molweight(ND2)
```

ans =

C: 1818
H: 2882
N: 420
O: 471
S: 25

ans =

3.8960e+004

For further investigation of the properties of the ND2 protein, try using **proteinplot**. This is a graphical user interface (GUI) that allows you to easily create plots of various properties, such as hydrophobicity, of a protein sequence. Click on the "Help" menu in the GUI for more information on how to use the tool.

```
proteinplot(ND2)
```

3. ANÁLISIS FILOGENÉTICO

- Proceso que usamos para determinar la relación evolutiva entre organismos.
- Estos resultados pueden dibujarse en un diagrama jerárquico llamado filograma (árbol filogenético)

EJEMPLO: CONSTRUCCIÓN DE UN ARBOL FILOGENÉTICO

- Usamos datos de secuencias mitocondriales (D-loop), creamos un árbol filogenético para una familia de primates.
- A partir de secuencias mitocondriales DNA (mtDNA) para la familia HOMINIDAE.

3.1. BÚSQUEDA DE DATOS FILOGENÉTICOS EN NCBI

PROCEDIMIENTO

- Use the MATLAB Help browser to search for data on the Web. In the MATLAB Command Window, type web('http://www.ncbi.nlm.nih.gov')
- Buscar el sitio web NCBI para información
- Seleccionar el link taxonomy para la familia HOMINIDAE

3.2. CREACIÓN DE UN ARBOL FILOGENÉTICO PARA CINCO ESPECIES

- Crear una estructura MATLAB con información acerca de las secuencias.
- Se emplean códigos de acceso para las secuencias mitocondriales D-loop aisladas de especies homínidas diferentes.

```
data = {'German_Neanderthal'      'AF011222';
        'Russian_Neanderthal'    'AF254446';
        'European_Human'        'X90314' ;
```

<http://matlab.udes.edu.co>

http://es.geocities.com/matlab_colombia/diamatlabnov3.html

<http://www.compelect.com.co/FormularioDiaMATLAB.html>

PONENCIA # 4

```
'Mountain_Gorilla_Rwanda' 'AF089820';
'Chimp_Troglodytes'      'AF176766'; };
```

- Obtener secuencias de datos de la base de datos GenBank y copiarlas en MATLAB

```
for ind = 1:5
    seqs(ind).Header = data{ind,1};
    seqs(ind).Sequence = getgenbank(data{ind,2},
                                    'sequenceonly', true);
end
```

- Se calculan las distancias de parejas discretas y se crea un objeto **phytree**

```
distances = seqpdist(seqs,'Method','Jukes-Cantor','Alphabet','DNA');
tree = seqlinkage(distances,'UPGMA',seqs)
```

- MATLAB dibuja un árbol filogenético

```
h = plot(tree,'orient','bottom');
ylabel('Evolutionary distance')
set(h.terminalNodeLabels,'Rotation',-45)
```

3.3. CREACION DE UN ARBOL FILOGENÉTICO PARA DOCE ESPECIES

- Se adicionan más secuencias a la estructura MATLAB

```
data2 = {'Puti_Orangutan'      'AF451972';
         'Jari_Orangutan'      'AF451964';
         'Western_Lowland_Gorilla' 'AY079510';
         'Eastern_Lowland_Gorilla' 'AF050738';
         'Chimp_Schweinfurthii'   'AF176722';
         'Chimp_Vellerosus'       'AF315498';
```

<http://matlab.udes.edu.co>
http://es.geocities.com/matlab_colombia/diamatlabnov3.html
<http://www.compelect.com.co/FormularioDiaMATLAB.html>

```
'Chimp_Verus'      'AF176731';  
};
```

- Obtener las secuencias adicionales de datos adicionales de la base de datos GenBank.

```
for ind = 1:7  
    seqs(ind+5).Header = data2{ind,1};  
    seqs(ind+5).Sequence = getgenbank(data2{ind,2},  
                                         'sequenceonly', true);  
end
```

- Se calculan las distancias pares iguales y el apareamiento jerárquico

```
distances = seqpdist(seqs,'Method','Jukes-Cantor','Alpha','DNA');  
tree = seqlinkage(distances,'UPGMA',seqs);
```

- MATLAB dibuja un árbol filogenético

```
h = plot(tree,'orient','bottom');  
ylabel('Evolutionary distance')  
set(h.terminalNodeLabels,'Rotation',-45)
```

4. Referencias:

- 4.1. Bioinformatics Toolbox For Use with MATLAB®, User Guide, V. 21.1, The MathWorks Inc. 2005.
- 4.2. MATLAB 7.1, Release 14 Service Pack 3, The MathWorks Inc. 2005.
- 4.3. Bioinformatics Toolbox 2.1.1. The MathWorks Inc. , Septiembre 2005,