

Bioinformática

Análisis Bioinformático de la Major urinary protein 1

VÍCTOR RODRÍGUEZ PASTOR¹

Universidad de Salamanca

victorrrp@usal.es

SUMARIO

El análisis bioinformático de las distintas proteínas es una técnica que se usa diariamente en muchos laboratorios de investigación. Recientemente, el número de proteínas conocidas ha aumentado exponencialmente gracias a los numerosos avances en los campos de la secuenciación masiva y la genómica funcional. Este trabajo se ha centrado en la “Major urinary protein 1”, una proteína realmente poco conocida y que por tanto presenta un mayor interés de ser analizada en cuanto a su evolución filogenética entre las distintas especies.

Palabras clave: Bioinformática, árbol filogenético, Major urinary Protein 1, evolución.

SUMMARY

The analysis of bioinformatic proteins is employed daily in most research laboratories. Due to the amount of massive sequencing and breakthroughs in functional genomics, the number of known proteins is rising at an exponential rate. The major urinary protein 1 was chosen due to lack of information with regards to its function and phylogenetic relations. Therefore, this study will explore its origins and chronology.

Key words: Bioinformatics, phylogenetic tree, Major urinary Protein1, evolution.

¹ Víctor RODRÍGUEZ PASTOR es estudiante de 3º del Grado en Biotecnología en la Universidad de Salamanca.

1. INTRODUCCIÓN

El análisis bioinformático realizado tuvo como objetivo elegir una proteína de interés, la “Mayor urinary protein 1”, e intentar averiguar sus relaciones filogenéticas en las distintas especies. Para ello se realizaron los siguientes análisis bioinformáticos: En primer lugar, se realizó un BLAST para ver las secuencias parecidas a ésta que existen en las distintas especies, para ello se retocaron los distintos parámetros hasta encontrar una configuración óptima. En segundo lugar, se guardaron y anotaron las distintas secuencias pertenecientes a cada especie y se hizo un alineamiento de múltiples secuencias (MSA) utilizando tres alineadores distintos para poder contrastar los resultados. Una vez realizado el alineamiento de múltiples secuencias, se utilizó el programa MEGA7 mediante tres algoritmos diferentes para la creación de un árbol filogenético que nos mostrara el emparejamiento de las distintas especies según la evolución de esta proteína y los antecesores predichos por el programa, con su correspondiente significación estadística. Por último, se estableció un patrón y una secuencia consenso, de manera que, si en el futuro se quiere introducir otra proteína en la misma familia que la “Mayor urinary protein 1”, se podría comprobar si ésta se ajusta.

En cuanto a la función de la “Mayor urinary protein 1”, solo se sabe que se une a feromonas que son excretadas en la orina de los machos y que, por tanto, tiene relación con el comportamiento sexual de las hembras. Además, es un receptor que se activa por insulina y éste está relacionado con muchas funciones biológicas como la respiración y la regulación del metabolismo.

El organismo en el que se ha identificado esta proteína es el *Mus musculus* (Ratón), su identificador en UniProt es P11588 y se abrevia como MUP 1.

2. MÉTODOS

2.1. BLAST

En primer lugar, hice un p-blast para tener una visión previa de las distintas especies que tenían proteínas con secuencias parecidas a la MUP 1. En los parámetros excluí todas aquellas proteínas provenientes del *Mus musculus* ya que sabemos que este es el organismo de procedencia de la proteína y queremos encontrar la relación entre distintas especies. Además, también se filtraron de la búsqueda todas aquellas secuencias que o bien solo hubieran sido predichas y no se hubiera obtenido en el laboratorio o bien se hubieran obtenido mediante metagenómica (*uncultured*).

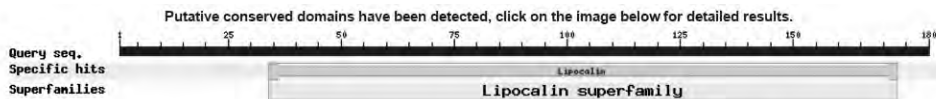


Figura 1. El BLAST nos informa de que tenemos un dominio perteneciente a la superfamilia de las lipocalinas.

Desde el primer momento, el BLAST ya predijo que la secuencia aminoacídica de la MUP1 contiene un dominio conservado de la superfamilia de las lipocalinas, a saber, varios miembros que comparten una estructura común, definida en general por la unión de ligandos lipofílicos.

Una vez realizado el BLAST, vemos que nos devuelve en primer lugar la misma proteína que no estaba bien anotada en la base de datos con el nombre de la especie y que por tanto no había sido filtrada al excluir *Mus musculus*. Una vez eliminada ésta, nos fijamos en que en el resto de secuencias encontradas, el porcentaje de identidad se encuentra entre el 40% y el 60%, a excepción de la perteneciente al organismo *Mus Spretus* que es otro tipo de ratón y por tanto también eliminaremos de la búsqueda.

El siguiente paso fue modificar los parámetros por defecto del BLAST ya que las identidades no son demasiado altas y por tanto podría haber secuencias que estuviéramos perdiendo por el camino al tener los parámetros demasiado estrictos: El *Word size* lo mantenemos en 6 ya que parece que bajarlo solo serviría para que el BLAST tardara más en realizarse, puesto que estamos encontrando un número suficiente de secuencias con este nivel de sembrado. El número máximo de secuencias lo subimos a 250 para poder tener mayor variedad de especies distintas y que nuestro árbol filogenético sea más rico. Además, desactivamos el botón de consultas cortas ya que no queremos que el BLAST nos modifique los parámetros. El *expected threshold* lo mantenemos en 10 ya que nos aporta un número suficiente de secuencias sin subirlo y la matriz la bajamos a una BLOSSUM más baja, en este caso, tras probar también varias PAM elegimos la BLOSSUM 45, para poder encontrar secuencias más divergentes que con la matriz por defecto. Por otro lado, mantenemos los valores por defecto de penalización de la creación y extensión de *gaps* ya que necesitaríamos tener un mayor conocimiento biológico sobre la proteína para poder retocar este parámetro de forma fiable.

Tras haber cambiado estos parámetros, nos encontramos que hay algunas secuencias nuevas que antes no estaban y además los e-valores de las que ya habían aparecido en la primera búsqueda son menores y por tanto estadísticamente más fiables (incluso hasta un valor de e-74).

Se recogen en la siguiente tabla las 36 secuencias elegidas para la realización del estudio bioinformático; si bien, de algunas especies cogimos más de una secuencia para posteriormente corroborar que los algoritmos de creación de árboles filogenéticos habían funcionado correctamente, ya que, *a priori*, las secuencias pertenecientes al mismo organismo deberían aparecer en ramas adyacentes.

| ESPECIE | Accession Number |
|---------------------------------|------------------|
| [Mus musculus] | P11588 |
| [Rattus norvegicus] | NP_671745.1 |
| [Rattus norvegicus] | NP_976070.1 |
| [Mesocricetus auratus] | AGV07608.1 |
| [Myotis davidii] | ELK24596.1 |
| [Felis catus] | NP_001009233.1 |
| [Canis lupus familiaris] | NP_001271390.1 |
| [Oryctolagus cuniculus] | CCC15303.1 |
| [Equus caballus] | NP_001075966.1 |
| [Bos taurus] | DAA26480.1 |
| [Bos mutus] | ELR50248.1 |
| [Heterocephalus glaber] | EHB13722.1 |
| [Sus scrofa] | NP_998979.1 |
| [Cavia porcellus] | CAX62131.1 |
| [Fukomys damarensis] | KFO26747.1 |
| [Sus scrofa] | 1GM6_A |
| [Capra hircus] | AHZ46504.1 |
| [Camelus ferus] | EPY89972.1 |
| [Cricetulus griseus] | ERE70560.1 |
| [Tupaia chinensis] | ELW70591.1 |
| [Macaca mulatta] | NP_001165312.1 |
| [Heterocephalus glaber] | EHB11625.1 |
| [Tachyglossus aculeatus] | ABJ90458.1 |
| [Cavia porcellus] | NP_001192113.1 |
| [Rattus norvegicus] | CAA29100.1 |
| [Oncorhynchus masou formosanus] | ABY21331.1 |
| [Fukomys damarensis] | KFO22774.1 |
| [Oreochromis niloticus] | CCF55069.1 |
| [Oncorhynchus mykiss] | CDQ80766.1 |
| [Anoplopoma fimbria] | ACQ59033.1 |
| [Salmo salar] | ACI69895.1 |
| [Oreochromis niloticus] | CCF55070.1 |

| | |
|-------------------|----------------|
| [Pteropus alecto] | ELK07359.1 |
| [Xenopus laevis] | NP_001081513.1 |
| [Danio rerio] | NP_998799.1 |
| [Cavia porcellus] | NP_001192113.1 |

Figura 2. Tabla con las secuencias escogidas.

Estos resultados se obtuvieron con los parámetros ya comentados anteriormente:

| Search Parameters | |
|-------------------------|----------|
| Program | blastp |
| Word size | 6 |
| Expect value | 10 |
| Hidist size | 250 |
| Gapcosts | 15,2 |
| Matrix | BLOSUM45 |
| Filter string | F |
| Genetic Code | 1 |
| Window Size | 40 |
| Threshold | 21 |
| Composition-based stats | 2 |

| Database | |
|---------------------|---|
| Posted date | May 6, 2016 8:16 PM |
| Number of letters | 31,916,365,256 |
| Number of sequences | 87,005,143 |
| Entrez query | all [filter] NOT[txid10090 [ORGN] OR txid10096 [ORGN] OR (XP_000001:XP_999999[paac] OR XP_00000001:XP_99999999[paac]) OR environmental samples[organism] OR metagenomes[orgn] OR txid32644[orgn]] |

| Karlin-Altschul statistics | | |
|----------------------------|-----------|---------|
| Lambda | 0.231236 | 0.203 |
| K | 0.0940427 | 0.041 |
| H | 0.244528 | 0.12 |
| Alpha | 0.9113 | 1.7 |
| Alpha_v | 9.61106 | 44.8257 |
| Sigma | | 45.0604 |

Figura 3. Parámetros del Blast.

También probé a hacer un PSI-BLAST con 5 interacciones, pero la mayoría de las secuencias nuevas respecto al p-BLAST eran sólo *hipotéticas* y, como ya tenía suficientes secuencias y especies distintas, decidí no incluir ninguna de ellas.

2.2. ALINEAMIENTO DE MÚLTIPLES SECUENCIAS (MSA)

El alineamiento de múltiples secuencias es la técnica bioinformática que nos permite aplicar un algoritmo y conseguir alinear secuencias aminoacídicas de proteínas que tienen una cierta homología. De esta manera podemos encontrar qué zonas están más conservadas en la evolución, lo cual indica que son regiones muy importantes para la función de la proteína, ya que las zonas menos importantes tienen mayor número de variaciones en la secuencia aminoacídica entre las distintas especies. Además, el alineamiento de múltiples secuencias sirve como paso previo para encontrar posibles sitios catalíticos de las enzimas que sean susceptibles de sufrir mutaciones dirigidas para aumentar su efectividad.

Para el alineamiento de múltiples secuencias usé tres herramientas bioinformáticas, de las cuales dos de ellas fueron *online* (Clustal Omega y T-coffee) y la otra fue el software informático Jalview.

En primer lugar, utilicé Clustal Omega en la página web del EBI (Instituto Europeo de Bioinformática). Esta herramienta no nos da demasiadas opciones en cuanto a cambio de parámetros se refiere, probé a aumentar el número de interacciones del HMM y el número de interacciones del árbol guía para así conseguir resultados más fiables. Aunque las zonas conservadas finales se presentarán con Jalview, aquí expongo una de las zonas conservadas obtenida con Clustal Omega (nuestra proteína de partida está identificada como “MUP1-Musmusculus” para así poder diferenciarla del resto y ver su homología):

| | |
|---------------------------------|---------------------------|
| gi 147904028 ref NP_001081513.1 | KNCTTIITPTAD-GNLEVTATVPI |
| gi 47174758 ref NP_998799.1 | KMGTAMLVPTQE-GDLDLSYANLI |
| gi 229368106 gb ACQ59033.1 | KMGTAMLVPTA--GDLDLSYANLI |
| gi 642093356 emb CDQ80766.1 | KMGTSVMPLTAG-GDLDLTYTNLI |
| gi 162949438 gb ABY21331.1 | KMGTSVMPLTAE-GDLDISSAMRI |
| gi 929318442 ref XP_014037121.1 | KMGTSIMLPTAG-GDLDISSAMRI |
| gi 542253638 ref XP_005462661.1 | KSGTAVVEPTED-GSLKVAFFSFP |
| gi 348518740 ref XP_003446889.1 | KTGTAIKPTED-GGIELSFSLSL |
| gi 326937536 ref NP_001192113.1 | RAYFRQVDCTEGCDIISITFYTFI |
| gi 676266951 gb KFO22774.1 | RVYFRQLCCQDGCINISIRFYV-I |
| gi 116248639 gb ABJ90458.1 | RCYMSSIDPAWIN-ESIRFNFYV-I |
| gi 50978944 ref NP_001003189.1 | RVFIHMSA-KD-GNLHGDIIL-I |
| gi 351708706 gb EHB11625.1 | HLFIRTIELLN--SLLFHFHF-I |
| gi 431898989 gb ELK07359.1 | RVFISISIQSLDN-GNLRFQFVL-I |
| gi 537147902 gb ERE70560.1 | RLFMRNHLLN--GSLKFDLFI-I |
| gi 633267632 gb AHZ46504.1 | RIFIESIQVIED-SGLKLSFHF-I |
| gi 946600638 ref XP_014407339.1 | RVFVQSIESLEN-GGLRFSFHF-I |
| gi 284519695 ref NP_001165312.1 | RVFVRNIEHLN--GSLKFDFFI-I |
| gi 519113840 emb CAX62131.1 | RVFVESIEPVID-SALSFKFMA-I |
| gi 351710803 gb EHB13722.1 | RVFVEFYALKN--SSIFFKFI-I |
| gi 676271779 gb KFO26747.1 | RVFVESIQALKN--SSVFFKFI-I |
| gi 444730201 gb ELW70591.1 | RIFVEHIDLLN--SSLSVNFHT-I |
| MUP1-Musmusculus MUP1_MOUSE | RLFLEQIHVLEN--SLVLFKFI-I |
| gi 22219448 ref NP_671745.1 | RVFVQHIVLEN--SLAFKFI-I |
| gi 42627893 ref NP_976070.1 | RVFVQHIVLEN--SLGPKFI-I |
| gi 56246 emb CAA29100.1 | RVFMQHIDVLEN--S----- |
| gi 557943216 emb CCC15303.1 | RVFVEYIHWKN--SSLSFKFHT-I |
| gi 880959776 ref XP_012979903.1 | RVFVKSIHLFN--SSLAFKFI-I |
| gi 47523218 ref NP_998979.1 | RVFVEHIVLDN--SSLAFKFI-I |
| gi 21465464 pdb 1GM6 A | RVFVEHIVLDN--SSLAFKFI-I |
| gi 61820389 ref XP_590993.1 | RVFVEYIDVLEN--SSLLFKFI-I |
| gi 440898827 gb ELR50248.1 | RFFVEYSLEN--SSLFKFI-I |
| gi 126723762 ref NP_001075966.1 | RVFVDIVRALDN--SSLYAEYQT-I |
| gi 432091571 gb ELK24596.1 | RVFVESIQVLYN--SSLSFKFHI-I |
| gi 57163775 ref NP_001009233.1 | RVFVEHIVRALDN--SSLSFKFI-I |
| gi 548923872 ref NP_001271390.1 | RVFVKDIEVLSN--SSLIFTMHT-I |

Figura 4. Alineamiento con Clustal omega.

El siguiente método elegido fue la herramienta T-coffee. A pesar de que este programa *online* tampoco nos permite cambiar parámetros interesantes, sí que nos decidimos a usar su versión PSI-Coffee. La elección de esta versión del programa se debe a que las 36 secuencias que hemos elegido para nuestro estudio son de especies muy diversas y, además, sus porcentajes de identidad ronda en algunos casos el 30%. Por tanto, es mejor utilizar esta versión que, a pesar de ser más lenta (tardó en su ejecución 4 horas y 59 minutos mientras que el T-Coffe normal se ejecutó en 10 segundos), es mucho más precisa a la hora de encontrar homologías en proteínas que a priori están distanciadas. Aquí ejemplifico una de las zonas conservadas encontradas:



Figura 5. Alineamiento con Psi-Coffee.

En último lugar, utilicé el programa Jalview para visualizar el alineamiento. Jalview es un software que te permite utilizar distintos algoritmos alojados en servicios web, tales como Clustal Omega, T-coffee, etc. En este caso, como ya habíamos utilizado los dos anteriores en sus versiones *online*, utilicé el algoritmo Muscle. Como en el BLAST, decidí cambiar la matriz a una BLOSSUM45 ya que estas secuencias tienen relativamente poca identidad y por tanto se obtienen mejores resultados si se es más permisivo con las sustituciones de aminoácidos. Además le indiqué que se trataban de secuencias proteicas. Sin embargo, el programa me devolvió errores de ejecución y por tanto probé con las matrices BLOSSUM40, PAM170 y PAM140 (matrices destinadas a secuencias lejanas) pero todas daban error. Parece que la última versión del jalview debe tener algún tipo de error con el Muscle ya que solo me devuelve resultados si ejecuto el algoritmo por defecto, lo cual es una pena porque el alineamiento podría ser más preciso con una matriz más permisiva. Por tanto, ejecuté el Muscle por defecto y obtuve las siguientes zonas conservadas:

En último lugar, utilicé el programa Jalview para visualizar el alineamiento. Jalview es un software que te permite utilizar distintos algoritmos alojados en servicios web, tales como Clustal Omega, T-coffee, etc. En este caso, como ya habíamos utilizado los dos anteriores en sus versiones *online*, utilicé el algoritmo Muscle. Como en el BLAST, decidí cambiar la matriz a una BLOSSUM45 ya que estas secuencias tienen relativamente poca identidad y por tanto se obtienen mejores resultados si se es más permisivo con las sustituciones de aminoácidos. Además le indiqué que se trataban de secuencias proteicas. Sin embargo, el programa me devolvió errores de ejecución y por tanto probé con las matrices BLOSSUM40, PAM170 y PAM140 (matrices destinadas a secuencias lejanas) pero todas daban error. Parece que la última versión del jalview debe tener algún tipo de error con el Muscle ya que solo me devuelve resultados si ejecuto el algoritmo por defecto, lo cual es una pena porque el alineamiento podría ser más preciso con una matriz más permisiva. Por tanto, ejecuté el Muscle por defecto y obtuve las siguientes zonas conservadas:

```

MUP1-Musmusculus[MUP1_MOUSE/32-60] I N G E W H T I I L A E K K E K I E W H N F L F L
gi|22219446|ref|NP_671745.1|/33-61 L N G D W F S I V V A S K K K I E H S M V F V
gi|42627893|ref|NP_976070.1|/33-61 L N G D W F S I M A A S K K K I E H S M V F V
gi|1800959776|ref|XP_012979903.1|/36-64 I V E C H S I L L A S Q K E M I E H S M V F V
gi|432091571|gb|ELK24596.1|/32-60 I S G E W Y S I L L A S D M E M I E H S M V F V
gi|57163775|ref|NP_001009233.1|/30-59 I S G E W Y S I L L A S V E K I E H S M V F V
gi|548923672|ref|NP_001271390.1|/31-59 I S G D W Y S I L L A S D I E K I E H S M V F V
gi|557943216|emb|CC015303.1|/19-47 D I S G E W S V L L A S D H E K I E H S M V F V
gi|126723762|ref|NP_001075966.1|/32-60 I S G E W Y S I F L A S D V E K I E H S M V F V
gi|161820389|ref|XP_590993.1|/31-59 E I A G E W Y S I L L A S N E K I E H S M V F V
gi|440898827|gb|ELR50248.1|/30-56 I T G E W F S I L L A S D N E K I E H S M V F V
gi|351710803|gb|EHB13722.1|/31-59 I S G K W Y S V L L A S K K E K I E D E S M V F V
gi|47523218|ref|NP_998979.1|/33-61 I A G E W Y S I L L A S A A E N I E H S M V F V
gi|519113840|emb|GAX62131.1|/31-59 I S G N W Y T V K E A S K K S T I E G G S M V F V
gi|676271779|gb|KFO26747.1|/31-59 I S G K W Y S V L L A S Q K K T I E G G S M V F V
gi|21465464|odb|1GM6|A/17-45 I A E E W S I L L A S A A E N I E H S M V F V
gi|633267632|gb|AH246504.1|/32-60 V S G T W Y S I S M A A N M K R I E E D D L I F I
gi|946600638|ref|XP_014407339.1|/32-60 V S G T W Y S I S M A T D M K R I E E D D L V F I
gi|537147902|gb|ERE70560.1|/32-60 I S R N W Y T I C M A S N M T R I E E N D L F V F
gi|444730201|gb|ELW70591.1|/31-59 I S G S W Y S I L T A S D N E K I G E S F R M I F M
gi|284519689|ref|NP_001165312.1|/33-31 I S G V W Y S I F M A S D L N R I K E N D L V F V
gi|351708706|gb|EHB11625.1|/32-60 P A A G S W S G I S L A S N V M W I G V N D L H F I
gi|116246639|gb|ABJ90458.1|/29-57 I S G R W I T L W M A A D T S L V M T H P L G V M
gi|50978944|ref|NP_001003189.1|/32-60 E L S G R W H S V A L A E N K S L I K R W H F V F I
gi|47174756|ref|NP_998799.1|/32-60 I V E R K W L L V G F A T N A Q W F V S H K D M M G T
gi|326937536|ref|NP_001192113.1|/23-51 I V P G N W R T A I A A H V K I I V N E L H A T I
gi|56246|emb|CAA29100.1|/33-61 L N G D W F S I V V A S K K K I E H S M V F V
gi|162949438|gb|ABY21331.1|/30-58 I M A E K K W A V G F A T N A Q W F M N R K A M M G T
gi|676266951|gb|KFO22774.1|/12-40 D I G D W R T H A V A A N V K I E G G E L V F I
gi|542253638|ref|XP_005462661.1|/32-606 T A G K W L L T G V C E N S W F V C R K A S I S G T
gi|642093356|emb|CQ060766.1|/30-58 I M A G R W I I V G F A T N A H W F V S H K A D L M G T
gi|229368106|gb|ACQ59033.1|/32-60 I M A E K K W I I V G F A T N A Q W F V N N K A M M G T
gi|929318442|ref|XP_014037121.1|/30-58 I M A G K W A V G F A T N A Q W F V K R K G G M M G T
gi|348518740|ref|XP_003446889.1|/32-60A I M A K K W L I G I G E N A Q W F V S R K A N M T G T
gi|431898989|gb|ELK07359.1|/1-19 - - - - - M A E D M R R I E E D D L V F I
gi|147904028|ref|NP_001081513.1|/33-61 V L G K K W G I G L A E N S N W F K R K S H M M A C T

```

Figura 6. Zona conservada número 1.

```

MUP1-Musmusculus[MUP1_MOUSE/116-1] L D E K H G I L R E N I D L L
gi|22219446|ref|NP_671745.1|/119-134 L D V A H G I T R N I D L T
gi|42627893|ref|NP_976070.1|/119-134 L D V A H G I T R N I D L M T
gi|1800959776|ref|XP_012979903.1|/123-1 L D K E N G I V K N I L D L T
gi|432091571|gb|ELK24596.1|/119-134 L D K K N G I T K E N I L D L T
gi|57163775|ref|NP_001009233.1|/116-123 Q Q E H G I V - N I L D L T
gi|548923672|ref|NP_001271390.1|/118-130 Q Q G M E I P K N I L D L T
gi|557943216|emb|CC015303.1|/106-121 L Q Q E R G I V E R N I L D L T
gi|126723762|ref|NP_001075966.1|/119-131 V Q K R G I V K E N I D L T
gi|161820389|ref|XP_590993.1|/118-133 I Q Q K Y G V V K N V I D L T
gi|440898827|gb|ELR50248.1|/117-132 I Q Q K Y G V V K N V I D L T
gi|351710803|gb|EHB13722.1|/150-165 F Q Q E N Q V V R C N I L D L T
gi|47523218|ref|NP_998979.1|/119-134 I Q Q Y G I I K E N I D L T
gi|519113840|emb|GAX62131.1|/117-132 F Q H K N G I G E A I I D M T
gi|676271779|gb|KFO26747.1|/117-131 F S Q K N G W V - S T V L D L T
gi|21465464|odb|1GM6|A/103-118 I D Q Q Y G I I K E N I D L T
gi|633267632|gb|AH246504.1|/119-134 A D K S H G L G P E K I I R F E
gi|946600638|ref|XP_014407339.1|/119-131 D K T Y S L G P E N I V T M S
gi|537147902|gb|ERE70560.1|/119-134 I Q Q K Y G L D S K I I N L T
gi|444730201|gb|ELW70591.1|/115-130 L D E K Y G I F E N V V D L T
gi|284519689|ref|NP_001165312.1|/90-100 T Q K K Y G L G P D N I V D L T
gi|351708706|gb|EHB11625.1|/119-134 A R R K Y G M G P R N T I N L A
gi|116246639|gb|ABJ90458.1|/115-130 F N T L W D I V K E N I T V M Q
gi|50978944|ref|NP_001003189.1|/118-130 C E D I D L H K I Q I V V L G
gi|47174756|ref|NP_998799.1|/123-138 F O L D T O I L R N I V M L R
gi|326937536|ref|NP_001192113.1|/111-119 V Q N A Q I P L E N I R Y V I
gi|56246|emb|CAA29100.1|/71-70 - - - - -
gi|162949438|gb|ABY21331.1|/121-136 F L D T O I L S D N I V F L R
gi|676266951|gb|KFO22774.1|/100-115 V A E E Y G I R V E N T R N I
gi|542253638|ref|XP_005462661.1|/121-130 L Q L Q S O I L R E N I V Y L F
gi|642093356|emb|CQ060766.1|/121-136 F L D T O I L S D N I A I L P
gi|229368106|gb|ACQ59033.1|/122-137 L L E T O I L R N I A I L P
gi|929318442|ref|XP_014037121.1|/121-130 F L D T O I L S D N I V L R
gi|348518740|ref|XP_003446889.1|/123-130 F L Q L T S I L R N I V L R
gi|431898989|gb|ELK07359.1|/78-90 T C I S P Y P T E S P T L - - -
gi|147904028|ref|NP_001081513.1|/123-130 F E A K S Q Q L A D E R L I L R

```

Figura 7. Zona conservada número 2.

| | |
|---------------------------------------|-----------------|
| MUP1-Musmusculus[MUP1_MOUSE/3-15 | MLLLLLLCLLTFLVC |
| gi 22219448 ref NP_671745.1 /3-15 | LLLLLLLCLGLTFLV |
| gi 42627893 ref NP_976070.1 /3-15 | LLLLLLLCLGLTFLV |
| gi 880959776 ref XP_012979903.1 /3-15 | LLLLLLVLVQLLELT |
| gi 432091571 gb ELK24596.1 /3-15 | LLLLLQGLTFLVCA |
| gi 57163775 ref NP_001009233.1 /3-15 | LLLLLQGLTFLVCA |
| gi 549923872 ref NP_001271390.1 /3-15 | LLLLLQGLTFLVMA |
| gi 557943216 emb CGG15303.1 /3-15 | LCAPKRGDAHS9S |
| gi 126723762 ref NP_001075966.1 /3-15 | LLLLLQGLTFLVCA |
| gi 61820389 ref XP_590993.1 /3-15 | LLLLLQGLTFLVCA |
| gi 440898827 gb ELR50248.1 /3-15 | LLLLLQLTLVQAO |
| gi 351710803 gb EHB13722.1 /3-15 | LLLLLQGLTFLFT |
| gi 47523218 ref NP_998979.1 /3-15 | LLLLLQLGLTFLP9 |
| gi 519113840 emb CAX62131.1 /3-15 | LQLLQLGLTFLCT |
| gi 676271779 gb KFO26747.1 /3-15 | LLLLLQGLTFLFT |
| gi 21465464 pad 1GMS /3-15 | EAGQDVVRSNFD9A |
| gi 633267632 gb AHZ46504.1 /3-15 | LLLLLQGLTFLVCA |
| gi 946600638 ref XP_014407339.1 /3-15 | LLLLLQGLTFLVFA |
| gi 537147902 gb ERE70560.1 /3-15 | VLPLVLLVLLAA |
| gi 444730201 gb ELLW70591.1 /3-15 | LLLLLQGLTFLVCA |
| gi 284519695 ref NP_001165312.1 /3-15 | RISGVWYSIFMA9S |
| gi 351708708 gb EHB11625.1 /3-15 | LLLLLQGLTFLVLA |
| gi 116248639 gb ABJ90458.1 /3-15 | TLLLGITLALVMA |
| gi 50978944 ref NP_001003189.1 /3-15 | LLLLLVGLTFLJCG |
| gi 47174758 ref NP_998799.1 /3-15 | GVVVKMLCLLCA |
| gi 326937536 ref NP_001192113.1 /3-15 | QLLLLALVSLAD |
| gi 56246 emb CAA29100.1 /3-15 | LLLLLQGLTFLV |
| gi 162949438 gb ABY21331.1 /3-15 | LSIMGVLLCATLA |
| gi 676266951 gb KFO22774.1 /3-15 | NLPSVNLFLQIDG |
| gi 542253638 ref XP_005462661.1 /3-15 | WRALLGVVCECL |
| gi 642093356 emb CDQ80766.1 /3-15 | LRMMGLLCAALV |
| gi 229368106 gb ACQ59033.1 /3-15 | NLLRMLGALMCOV |
| gi 929318442 ref XP_014037121.1 /3-15 | LSIMGVLLCTFLV |
| gi 348518740 ref XP_003446889.1 /3-15 | SLHLLGVVLCGL |
| gi 431898989 gb ELK07359.1 /3-15 | SDDMRRIEEDGDL |
| gi 147904028 ref NP_001081513.1 /3-15 | RLLLALLSVAAGD |

Figura 8. Zona Conservada número 3.

| | |
|--|------------|
| MUP1-Musmusculus[MUP1_MOUSE/115-120 | DNFLMATHL |
| gi 22219448 ref NP_671745.1 /116-124 | NRRVMTIHL |
| gi 42627893 ref NP_976070.1 /116-124 | NRRVMTIHL |
| gi 880959776 ref XP_012979903.1 /120-125 | NEYIITIKL |
| gi 432091571 gb ELK24596.1 /116-124 | NVDLILLTL |
| gi 57163775 ref NP_001009233.1 /113-120 | NEFIILLHL |
| gi 549923872 ref NP_001271390.1 /115-120 | NEDFITIHL |
| gi 557943216 emb CGG15303.1 /103-111 | FNDYIIFHL |
| gi 126723762 ref NP_001075966.1 /116-120 | NEHILLLTL |
| gi 61820389 ref XP_590993.1 /115-123 | SQHIIFHL |
| gi 440898827 gb ELR50248.1 /114-122 | SQHIIFHL |
| gi 351710803 gb EHB13722.1 /113-121 | NMYIIFHL |
| gi 47523218 ref NP_998979.1 /116-124 | SDYVILLHL |
| gi 519113840 emb CAX62131.1 /114-122 | NKNVAFQQL |
| gi 676271779 gb KFO26747.1 /114-122 | PSSYVIFQQL |
| gi 21465464 pad 1GMS /100-108 | NSDVFILLHL |
| gi 633267632 gb AHZ46504.1 /116-124 | NQRVVIILHM |
| gi 946600638 ref XP_014407339.1 /116-120 | NRLIITFHL |
| gi 537147902 gb ERE70560.1 /116-124 | NRTVVFILHM |
| gi 444730201 gb ELLW70591.1 /112-120 | NKQFVQLCE |
| gi 284519695 ref NP_001165312.1 /87-95 | NRLFITFHL |
| gi 351708708 gb EHB11625.1 /116-124 | NWMTIIFVL |
| gi 116248639 gb ABJ90458.1 /113-120 | NNDLMLHT |
| gi 50978944 ref NP_001003189.1 /115-123 | NKSLILYFM |
| gi 47174758 ref NP_998799.1 /120-128 | NDEKAFHT |
| gi 326937536 ref NP_001192113.1 /108-116 | NMAFVIGVM |
| gi 56246 emb CAA29100.1 /71-70 | |
| gi 162949438 gb ABY21331.1 /116-126 | NDFLLIHT |
| gi 676266951 gb KFO22774.1 /97-105 | QNVILLNIT |
| gi 542253638 ref XP_005462661.1 /118-126 | NDFLNIHY |
| gi 642093356 emb CDQ80766.1 /118-126 | NDFALIHT |
| gi 229368106 gb ACQ59033.1 /119-127 | NQDVALVHT |
| gi 929318442 ref XP_014037121.1 /118-126 | NDFALIHT |
| gi 348518740 ref XP_003446889.1 /120-128 | DEHGLTHT |
| gi 431898989 gb ELK07359.1 /75-83 | NFLVITHT |
| gi 147904028 ref NP_001081513.1 /120-128 | NEVILMHT |

Figura 9. Zona conservada número 4.

He de comentar que los resultados obtenidos con los distintos algoritmos no variaron demasiado, las zonas conservadas quedaron prácticamente iguales y solo hubo cambios en las zonas de gaps. Esto se debe a que los porcentajes de identidad, aunque eran bajos (entre el 30 y el 40%), sí que se obtiene un mejor alineamiento que si se hubieran cogido secuencias con porcentajes de identidad cercanos al 20% que hubieran estado en el umbral de incertidumbre de si esas secuencias hubieran tenido ancestros comunes o no, y por tanto, los algoritmos informáticos no hubieran tenido un resultado claro y hubieran variado más.

En último lugar, se utilizó la herramienta WebLogo para ver las secuencias consenso en cada una de las cuatro zonas conservadas obtenidas anteriormente:



Figura 10. Logo zona conservada número 1.



Figura 11. Logo zona conservada número 2.

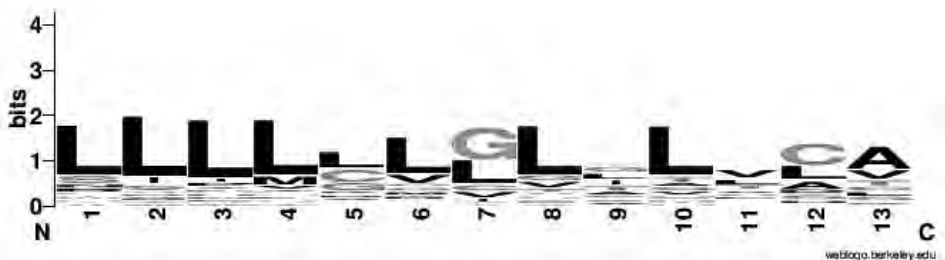


Figura 12. Logo zona conservada número 3.

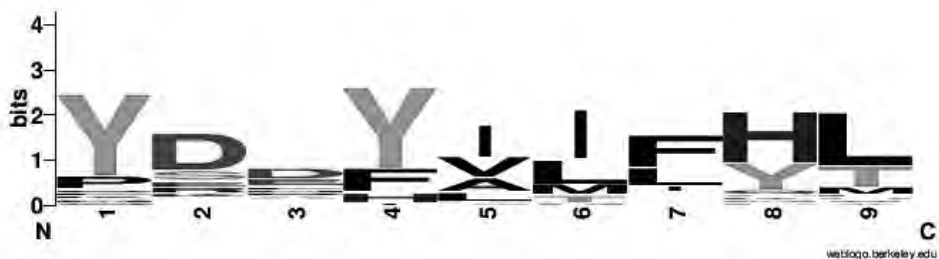


Figura 13. Logo zona conservada número 4.

La creación de estos Logos se basa en las proporciones relativas de cada aminoácido en cada columna del alineamiento de las regiones conservadas. De esta manera, las letras más grandes nos indican que hay una mayor proporción de ese aminoácido en esa columna del alineamiento de múltiples secuencias.

2.3. ÁRBOLES FILOGENÉTICOS

En este paso reside la mayor parte del interés del trabajo, ya que es donde por fin vemos las relaciones filogenéticas que existen entre las distintas especies teniendo en cuenta cómo han evolucionado estas secuencias de proteínas que tienen una cierta identidad entre sí.

Se han utilizado tres algoritmos distintos para la creación del árbol: NJ, UPGMA y máxima parsimonia. Todos ellos disponibles en el programa gratuito MEGA 7.

Además, le he aplicado a los tres el bootstrapping, un análisis estadístico que nos dice cómo de fiable es cada una de las uniones entre ramas del árbol.

Para realizar los árboles utilicé la zona conservada número 1 ya que es la más larga y, por tanto, la más significativa estadísticamente.

Al intentar hacer los árboles por primera vez se encontró un error en 2 de los 3 algoritmos (NJ y UPGMA) ya que sólo nos permite hacer árboles si todas las secuencias tienen aproximadamente la misma longitud, por tanto, tuve que desactivar la secuencia de *Pteropus alecto* para que me permitiera calcular correctamente las distancias. El problema de esta secuencia era que en esta zona conservada entre proteínas tenía 9 gaps seguidos al principio.

También conviene recordar que para leer las distancias existentes entre especies mediante el árbol filogenético, no se considera la distancia en vertical, sino en horizontal, ya que dos ramas pueden rotar sobre el nodo sin que cambie el signifi-

cado del árbol. La distancia se calcula sumando toda la distancia horizontal hasta el nodo común a la distancia existente hasta el otro miembro con el que queremos ver la distancia filogenética.

En primer lugar, utilicé el método NJ (Neighbor-Joining). Marqué la casilla del Bootstrap method, le informé de que se trataba de una secuencia de aminoácidos y fijé el número de interacciones del bootstrap en 500: (Ver ANEXO 1)

En este árbol podemos comprobar que el algoritmo se ha realizado correctamente ya que las secuencias pertenecientes a la especie *Sus Crofa* quedan muy próximas entre ellas, así como las de *Rattus Novergicus* y las de *Oreochromis niloticus*. Esto se debe a que, al pertenecer a la misma especie, han tenido menos tiempo evolutivo y por eso las secuencias son más parecidas. En el caso de *Rattus norvegicus*, a pesar de llamarse igual vemos que las secuencias son diferentes y por ello pensamos que han sido secuenciadas a partir de dos animales distintos pertenecientes a la misma especie. Por otro lado, vemos cómo las conservaciones a nivel de género también se mantienen a nivel de secuencia en los casos de *Oncorhynchus* y *Bos*. También es razonable que nuestra secuencia de partida de *Mus musculus* (es decir, ratón) se encuentre cercana a *Rattus novergicus* ya que los ratones y las ratas son muy parecidas.

Sin embargo, lo que más nos llama la atención son los casos de *Canis lupus familiaris* y de *Cavia porcellus* en los que a pesar de ser proteínas pertenecientes a la misma especie, están muy alejadas filogenéticamente. Esto se puede explicar porque se trata de alérgenos los cuáles varían bastante.

En siguiente lugar se realizó el método UPGMA también utilizando 500 interacciones de Bootstrapping y esta vez cambiamos el formato del árbol a straight: (VER ANEXO 2)

En este árbol vemos cómo se cumplen las mismas cosas que comentamos en el anterior. Esto confirma que los algoritmos se han llevado a cabo correctamente y aunque las significaciones estadísticas varíen, la predicción filogenética es buena.

Ahora en último lugar usamos el método de máxima parasimonia. Para la realización de este método sí que nos permite utilizar todas las secuencias y por tanto activamos otra vez la de *Pteropus alecto* y, además, como éste es el único en el que tenemos todas las secuencias, aumentamos el bootstrapping a 1000 para tener una mayor significación estadística: (VER ANEXO 3)

En este método lo que más se destaca es que *Mus musculus* aparece más alejado filogenéticamente de *Rattus novergicus* que en los casos anteriores y sobre todo el caso de las proteínas relativas a *Oncorhynchus* que aparecen muy alejadas entre ellas. Sin embargo, como este método tuvo un bootstrapping más grande que

los anteriores y tardó en realizarse casi 5 horas el árbol, entendemos que debería ser a priori más correcto que los dos anteriores. Aun así, sabemos que todo este tipo de predicciones filogenéticas no son dogmas, sino que se basan en predicciones estadísticas y por tanto los resultados varían entre métodos y además no asegura que la evolución ocurriera así realmente.

2.4. BÚSQUEDA DE PATRONES Y PHMMs

En este caso nos centramos en buscar un patrón para cada una de las zonas conservadas encontradas durante el alineamiento de múltiples secuencias para que así, en un futuro, si queremos ver si una nueva proteína se ajusta a esta familia de proteínas y por tanto tiene identidad con las proteínas de este grupo, solo tendremos que ver si se ajusta al patrón. Además, este patrón al introducirse en PROSITE nos deja ver qué proteínas relacionadas lo cumplen.

En cuanto a la primera secuencia conservada, como no todas las secuencias varían en las mismas posiciones y tienen porcentajes relativamente bajos de identidad entre ellas, fue difícil encontrar un patrón válido. En primer lugar, subí la *max flexibility* y el *max num flex spaces* a 5 cada uno para que hubiera más flexibilidad a la hora de hacer el patrón. Además, bajé el porcentaje mínimo de coincidencia con el patrón al 70%. El patrón resultante de estos cambios era demasiado corto y por tanto volví a bajar el porcentaje mínimo a 60% y nos dio lo siguiente:

| Best Patterns : | | | |
|-----------------|------------|------------|--|
| | fitness | hits(seqs) | Pattern |
| A | 1: 27.6904 | 22(22) | A-x-D-x(4)-I-x(0,1)-E-x(1,2)-G-x(2)-R-x(1,2)-F |
| B | 2: 27.6904 | 22(22) | A-x(1,2)-D-x(3,4)-I-E-x(2)-G-x(2)-R-x(1,2)-F |
| C | 3: 27.1904 | 22(22) | A-x-D-x(4)-I-x(0,1)-E-x(1,2)-G-x(2)-R-x(0,2)-F |

Figura 14. Patrones zona conservada número 1.

De los tres patrones con mayor significancia estadística cogí el primero y lo introduje en PROSITE para buscar proteínas que cumplieran este patrón. Nos salieron 16 secuencias encontradas, además, sabemos que el patrón funcionó porque alguna de estas secuencias encontradas en PROSITE estaban ya incluidas en nuestro análisis de múltiples secuencias como Major allergen Equ_c 1. *Equus caballus* (Horse). Además, también podemos confirmar que nos encontró proteínas relacionadas ya que nos salían en la búsqueda otras Mayor Urinary Proteíns del *Mus musculus* como la 11, la 6 y la 5.

En cuanto a la secuencia número 2, primero probé con un porcentaje mínimo del 80% pero me dio un patrón muy corto. Después probé con el 60% y me dio este resultado:

| Best Patterns : | | | |
|-----------------|---------|------------|---------------------|
| | fitness | hits(seqs) | Pattern |
| A 1: | 12.0102 | 24(24) | G-x(4)-N-x(0,1)-I |
| B 2: | 11.5102 | 23(23) | N-x(0,1)-I-x(1,2)-L |
| C 3: | 11.5102 | 25(22) | G-x(2,4)-N-I |

Figura 15. Patrones zona conservada número 2.

El problema de esta y de las dos secuencias conservadas que nos quedan por mostrar es que son más cortas que la primera y por tanto son menos específicas de esta familia de proteínas y el PROSITE no es capaz de encontrar específicamente este tipo de proteínas al tener patrones cortos y con poca puntuación.

De todas formas, mostraremos cuáles fueron los patrones de las secuencias conservadas número 3 y 4:

| Best Patterns : | | | |
|-----------------|---------|------------|-------------------------------------|
| | fitness | hits(seqs) | Pattern |
| A 1: | 18.8503 | 23(22) | L-x(2)-L-x(0,1)-G-x(0,2)-L-x(1,2)-L |
| B 2: | 18.8503 | 40(22) | L-x(2,3)-L-G-x(0,2)-L-x(1,2)-L |
| C 3: | 16.1802 | 44(23) | L-l-x(3)-l-x(1,2)-l |

Figura 16. Patrones zona conservada número 3.

| Best Patterns : | | | |
|-----------------|---------|------------|------------|
| | fitness | hits(seqs) | Pattern |
| A 1: | 7.8401 | 26(22) | Y-x(4,5)-L |
| B 2: | 7.8401 | 24(22) | Y-x(3,4)-L |
| C 3: | 7.3401 | 27(25) | Y-x(2,4)-L |

Figura 17. Patrones zona conservada número 4.

2.5. BÚSQUEDA EN PFAM DE LA SECUENCIA CONSENSO

Por último, descargamos la secuencia consenso del alineamiento con JALVIEW obtenida mediante el algoritmo MUSCLE y eliminamos toda la parte intermedia que solo existía en una proteína y por tanto no era representativa de la familia. Además, quitamos los signos “+” que no los reconoce como secuencia proteica el PFAM. La secuencia es la siguiente:

MMKLLLLLLCLGLLTLVCAHAEEAEDVVRSNFDLEKISGEWYSILLASDKK
 EKIEENGSMRVFVESIVLENS+SLSFKFAHTKVNGECTEVSIVCDKTEKDGV
 YTVEYQRWFDGENKFRIVETDYDDYIIFHLINFKNGETFQLLELYGRPDLSP
 ESLKEKQVQFCQKYGIVKENIIDLTKVDLRCLQARGSGVAQASSAETSD

Al introducirlo en PFAM, efectivamente nos dice (como ya nos había adelantado el BLAST al principio del trabajo) que nuestra secuencia consenso pertenece a la familia de las lipocalinas: “Lipocalin / cytosolic fatty-acid binding protein family”.

3. CONCLUSIONES

Nuestra investigación sobre la MUP1 (Major Urinary Protein 1) de *Mus musculus* se ha llevado a cabo correctamente. Se han variado los parámetros de los distintos algoritmos y procedimientos de forma óptima para obtener resultados interesantes sobre nuestra proteína que pertenece a la familia de las lipocalinas.

Se ha obtenido un árbol bastante rico con 35/36 secuencias (dependiendo del algoritmo usado) pertenecientes a 30 especies distintas en las que se refleja la evolución filogenética de proteínas parecidas a la nuestra de interés.

Además, se ha conseguido satisfactoriamente un patrón que encuentra otras proteínas parecidas en PROSITE y una secuencia consenso que se puede utilizar en el buscador de PFAM.

4. BIBLIOGRAFÍA

- ALTSCHUL, S.F., GISH, W., MILLER, W., MYERS, E.W. & LIPMAN, D.J. (1990) “Basic local alignment search tool”. *J. Mol. Biol.* 215: 403-410.
- ALTSCHUL, S.F., MADDEN, T.L., SCHÄFFER, A.A., ZHANG, J., ZHANG, Z., MILLER, W. & LIPMAN, D.J. (1997) “Gapped BLAST and PSI-BLAST: a new generation of protein database search programs”. *Nucleic Acids Res.* 25: 3389-3402.
- T-Coffee: A novel method for multiple sequence alignments. *Notredame, Higgins, Heringa, JMB*, 302 (205-217) 2000.
- EDGAR, ROBERT C. (2004), MUSCLE: multiple sequence alignment with high accuracy and high throughput, *Nucleic Acids Research* 32 (5), 1792-97.
- CROOKS GE, HON G, CHANDONIA JM, BRENNER SE. WebLogo: A sequence logo generator, *Genome Research*, 14:1188-1190, (2004).

- Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega Sievers F, WILM A, DINEEN DG, GIBSON TJ, KARPLUS K, LI W, LOPEZ R, MCWILLIAM H, REMMERT M, SÖDING J, THOMPSON JD, HIGGINS D. *Molecular Systems Biology* 7 Article number: 539 doi:10.1038/msb.2011.75
- MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets (Kumar, Stecher, and Tamura 2015).
- WATERHOUSE AM, PROCTER JB, MARTIN DMA, CLAMP M, BARTON GJ (2009) Jalview Version 2-a multiple sequence alignment editor and analysis workbench. *Bioinformatics* 25: 1189-1191. doi:10.1093/bioinformatics/btp033
- MICHENER, C.D., SOKAL, R.R. (1957): A quantitative approach to a problem of classification. *Evolution*, 11: 490-499.
- SAITOU, N., NEI, M. (1987). The neighbor-joining method: a new method for reconstructing phylogenetic trees, *Molecular Biology Evolution* 4: 406-425.
- “Performance of maximum parsimony and likelihood phylogenetics when evolution is heterogeneous”. *Nature* 431 (7011): 980-984
- SIGRIST CJA, DE CASTRO E, CERUTTI L, CUCHE BA, HULO N, BRIDGE A, BOUGUELERET L, XENARIOS I. New and continuing developments at PROSITE. *Nucleic Acids Res.* 2012; doi: 10.1093/nar/gks1067
- The Pfam protein families database: towards a more sustainable future: R.D. FINN, P. COGILL, R.Y. EBERHARDT, S.R. EDDY, J. MISTRY, A.L. MITCHELL, S.C. POTTER, M. PUNTA, M. QURESHI, A. SANGRADOR-VEGAS, G.A. SALAZAR, J. TATE, A. Bateman *Nucleic Acids Research* (2016) Database Issue 44:D279-D285.

5. ANEXOS

5.1. ANEXO 1: NEIGHBOR-JOINING

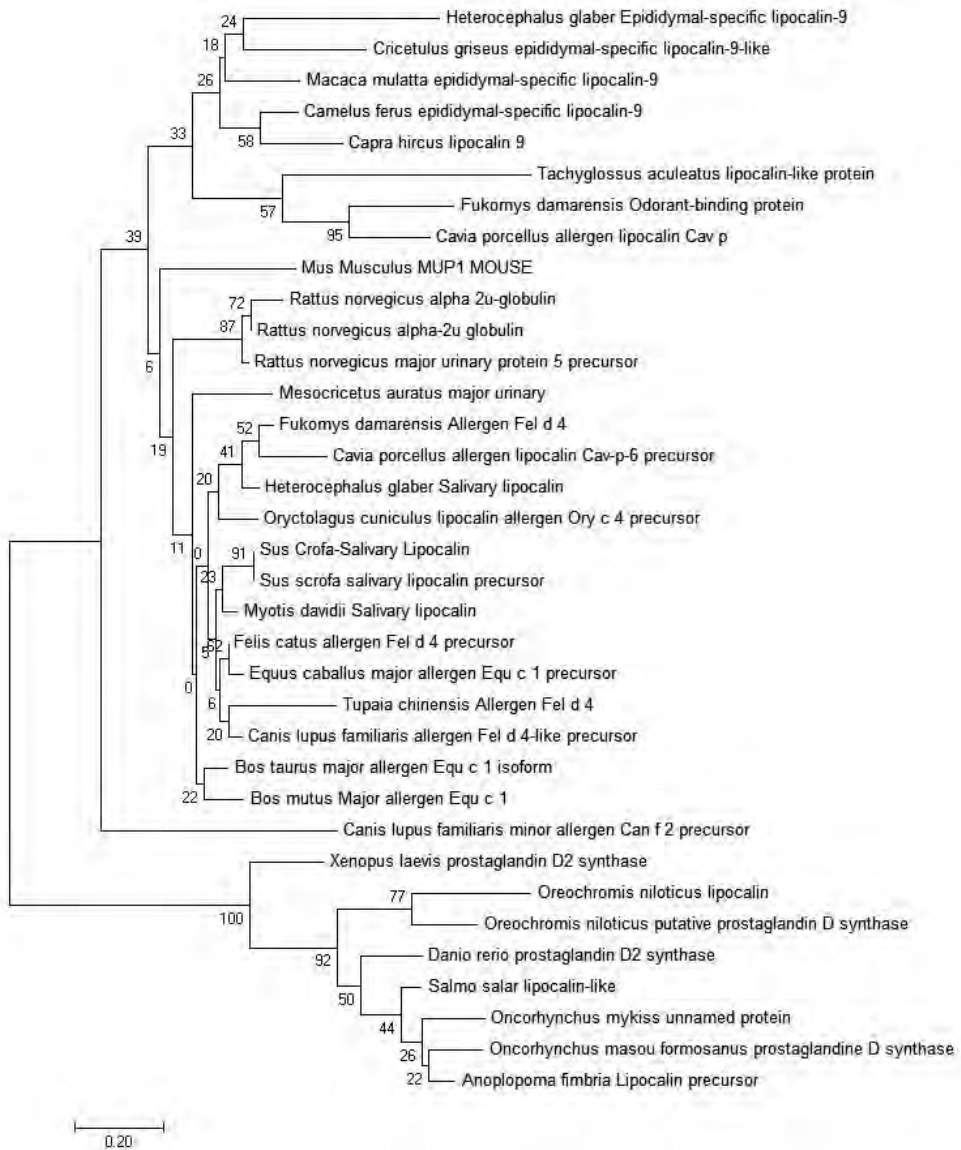


Figura 18. Árbol Neighbor-Joining.

5.2. ANEXO 2: UPGMA

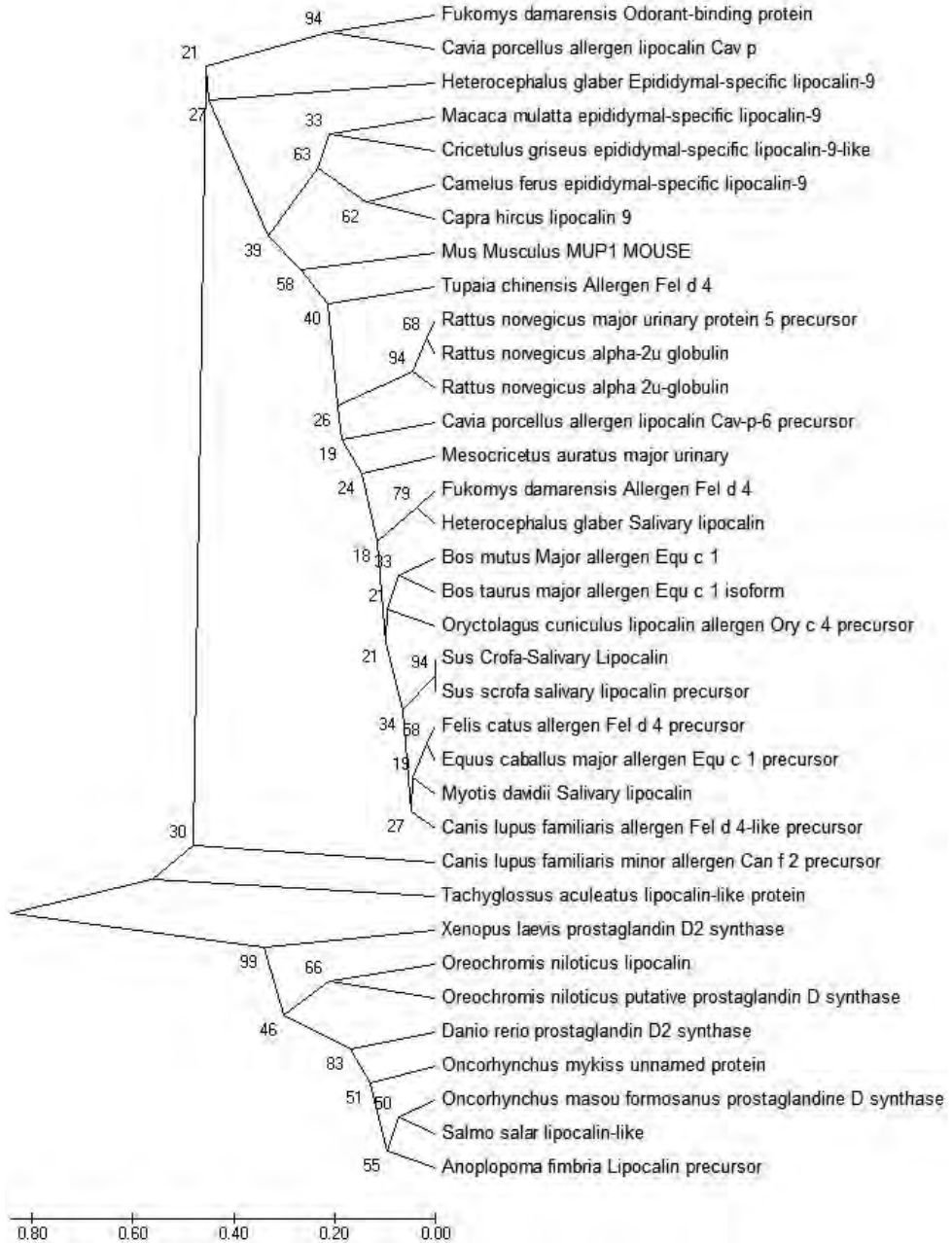


Figura 19. Árbol UPGMA.

5.3. ANEXO3: MÁXIMA PARASIMONIA

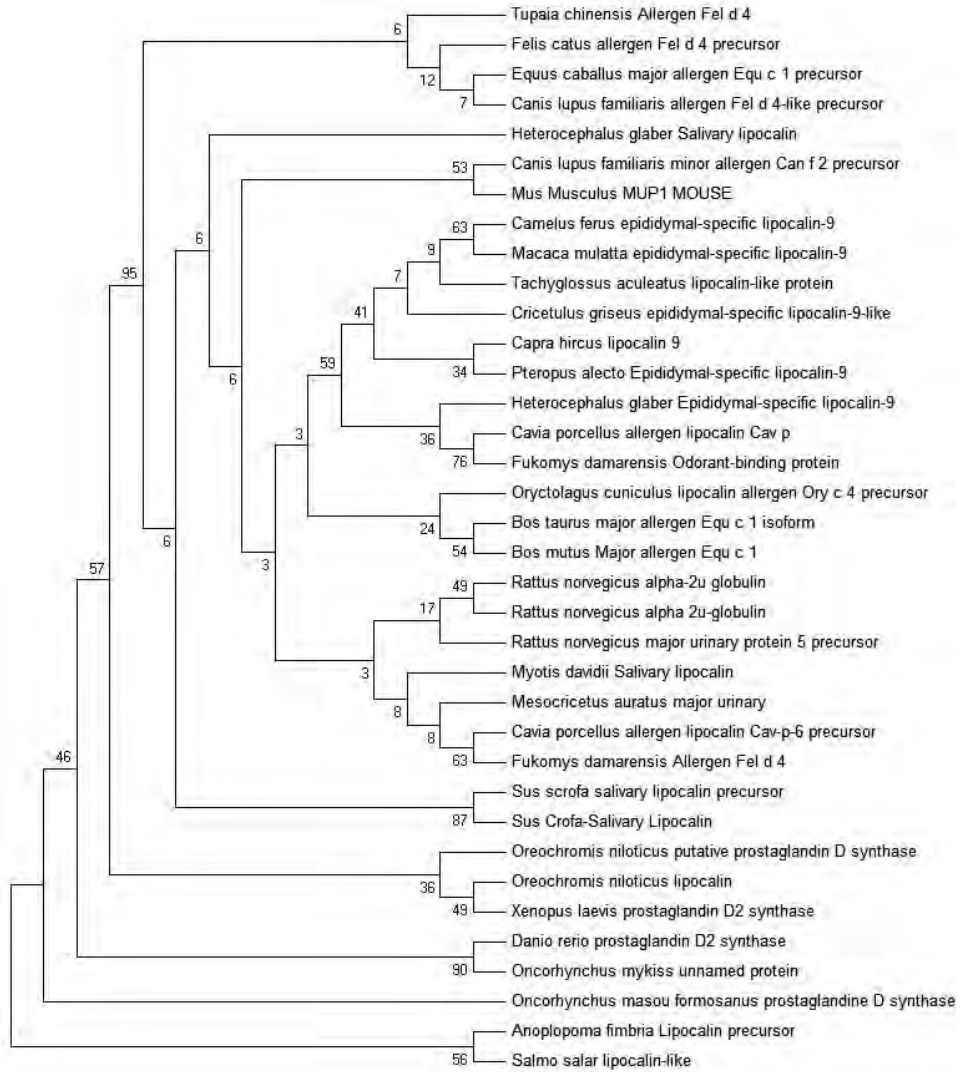


Figura 20. Árbol máxima parasimonia.