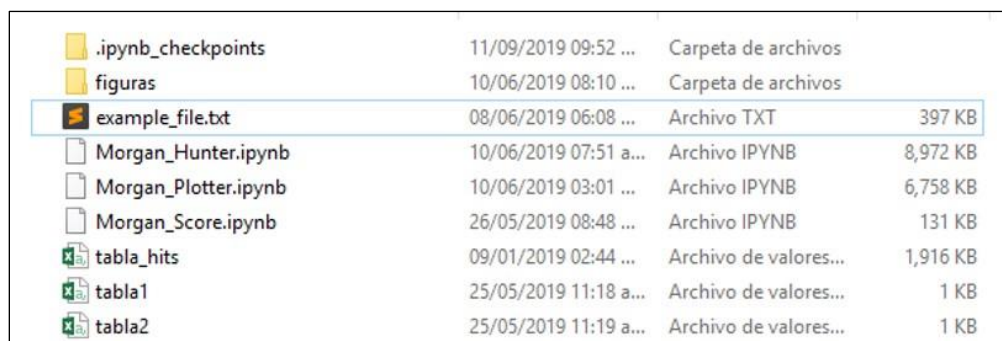# MORGAN HUNTER TOOL – USER MANUAL

The Morgan Hunter Tool is a Python-based computational tool that searches for motifs on DNA sequences using consensus sequences represented as Position Specific Frequency Matrices. The program provides some of the results in .csv files, others as graphical representations of the analyzed sequence(s), and yet others are displayed on screen.

The Morgan Hunter Tool is divided in three modules:

- Morgan Hunter (named after Sir Henry Morgan, the navigator and Dr. Thomas Hunt Morgan, the geneticist) searches for motif occurrences along DNA input sequences.
- Morgan Plotter produces the graphical representation of the analyzed sequences, with the occurrence of each motif found and the positions of satellite sequences (these must be provided by the user). The similarity of each detected motif with the consensus sequence is color coded. Morgan Plotter uses Morgan Hunter's output information, and therefore, these modules must be run sequentially.
- Morgan Score is used when the motif/consensus similarity of large amounts of sequences is to be found. It does not provide a graphical representation of the sequences. Morgan Score uses Morgan Hunter's output information, and therefore, these modules must be run sequentially.

## RUNNING THE MORGAN HUNTER TOOL

Download Anaconda Distribution (https://www.anaconda.com/distribution/). Make sure that the folder containing the three Morgan Hunter Tool modules (Morgan_Hunter.ipynb, Morgan_Plotter.ipynb, and Morgan_Score.ipynb) also contains the following items (Fig. 1):



*Figure 1*

- Morgan Hunter's input: a .txt file with the sequence(s) to be analyzed. The file layout must comply with the following format (example shown in Fig. 2):
    - sequence name preceded by ">" sign, line break
    - sequence without spaces, double line break.

```
>X07685.1
CATTCTCAGAAACTTCTTTGTGATGTGTGCATTCAACTCACAGAGTTGAACCTTCCTTTTCATAGAGCAG
TTTTGAAACACTCTTTTTGTAGAATCTGCAAGTGGATATTTGGACCGCTTTGAGGCCTTCGTTGGAAACG
GGAATATCTTCATATAAAAACTAGACAGAAG

>S67971.1
CATTCTCAGAAACTTCTTTGTGATGTGTCCATTCAACTCACAGAGTTGAACCTTTCTTTTGATAGAGCAG
TTTTGAAACACTCTTTTTGTAGAATCTGCAAGTGGATATTTGGAGCGCTTTGAGGCCTATGGTGGAAAAG
GAAATATCTTCACATAAAAACTAGACAGAAG

>AJ007752.1
GAATTCTCTAAAATTTCTTTCTGATGTGTGCATTCAACTCATAGTCTTAAACTTTTCTTTTGATAGAGCA
CTTTTGAAACATTCTTTTTGTAGAATCTGCAAGTGGATATTTGGAGCGCTTTGAGGCCTATGGTGGAAAA
GGAAATATCTTCACATAAAAGCTAGACAGAAGCATTCTCCGAAACTTCTTTGTGATATTTGCATTCAATT
CACCGATTTAAACTTTTCTTTTTATATAGCAGTTTTGAAACACTCGTTTAGTAGAATCTGCAAGTGGATA
TTTTGTTTCCTTTGAGGCCTATTTTGGAAAAGGAAATATCTTCACATAAAATATAACAGAAGCATTCTCA
GAAACTTGTTTTTGATGTGTGCATTCAACTCACGGAGATTAACGTTTCTTTTGACAGAGGAGTTTTGAAA
CACTCTTTTTGAAGAATCTGCAAGTGGATATTTGGTTTCCTTTGAGGCCTATGTTGGAAAATGAAATATC
TTCAAATAAAAACTAGACAGAAGCATTCTCACCAGCTTCGTTTTGATGTGAGCATTCAACTCCCAGAGTG
GAACCTTTCTTTTGATAGAGCAGTTTTGAAACACTCTTTTTGTGGAATCTGCAAGTGGATATTTGGAGAA
GTTTCAGGCCTATGGTGGAAAAGGAAATATCTTCACATAAAAGTAGAAGCATTCTCAGAAACTTCTTTG
TGATGTGTGCATTCAACACACAGAGTTGAACCTTTCTTCTGATAGAGCAGTTTTGAAACACTCTTTTTGT
AGTATCTGCAAGTGGATATTTGGAGCGCTTGGAGGCCTGTGGCGGAAAACGAAATATCTTCACATAAACA
CTAGACAGAAGAGTTCTCAGAAACTTCTTTGTGATGTGTGCATTCAACTCACAGAGTTCAACCTTTCTTC
TGATAGAGCAGTTCTGAGACACTCTTTTTGTAGAATTTGCAAGTGGGTATTTGGAGCACTTTGAGGCCTA
TTGTGAAAAGGAAATATCTTCGCATAAAAACTAGACAGAAGCATTCTCCAAAACTTCTTTGTGATGTGTG
CATTCAACTCACAGAGTTGAACCTTGCTTTTGATAGAGCAGTTTTGAAACACTCTTTCTGTAGAATCTGC
AAGTGGATATTTTGTTCCCTTTGAGGCCTATGGTGGAAAACGAAGTATCTTCACCCAAAAACTAGACAGA
AGCATTCTCAGAAACTTCT

>AK056737.1
AACATAGAGCTTCTTTCACTGGAGGACGAGAATAGGTGTGGCCTCAAAGTGCTCCAAATATCCACTTGCA
GATTCAGCAAAAAGAGTGTTTCAAAACTGCTCTATTAAAAGGAAGGTTCAGCTCTCTGTGTTGAATGCAC
ACATCACAAAGATGTTTCTGAGAATGCTTCTGTCTAGTTTTTATGTGAAGATATCCCCGTTTCCAACGAA
GTCCTCAGAGTGATCCAAATATCCACTTGCAGATTCTACAAAAAGAGTGATTCAAAACTGCTCTATCAAA
AGAAATGTTCAACTCTTTGAGTTGAATGCACATATCACAAAGCAGTTTCTGAGAATGCTTCTGTCTAGTG
TTTATGTGAAGATATTTCCTTTTCTACCACAGGCTTCAAATTACTCCAAATATCCACTTGCAAATTCTAC
AAAAAGAGTGTTTCATAACTGCTCTATCAAAAGGAAGGTTCAACCCTCTGTGTTGAATGCACACGTCGCA
AAGAAGTTTCTGAGAATGCTTCTGTCTAGTTTTTATGTGAAGATATTCCTGTTTCCACCGTAGGAATCAA
AGAGCTCCAAAGATCCACTTGCAGATTCTACAAAAAGAGGGCTTCAAAACTGCTCTATCAAAAGAAAGGT
TCAACTCTGTGAGTTGAATGCACACATCACCAAGCAGTTTCTGAGAATGCTTCTGTCTAGTTTTTTATATT
AAGATATTTCCTTTTCTACCATAGGCCTCAAAGCGCTCCAAATATTCACTTGCAGATTCTACAAAAAGAG
TGTTTCAAAACTGCTCTATGAAAAGGAAGTTTCAAATCTCTGTGTTGAATGCACACATCAAAAAGAAGTT
TCTGAGAATGCTTCTGTCTAGTTTTTATATGAAGATATTCCCATTTCCAACGAAGGCCTCAAAGCAGTCC
AAATATCCACCTGCAGATTCAACAAAAAGAGTGTTTCAAAACTGCTCTATCAAAAGAAAGGTTCAACTCT
GTGAGTTGAATGCACACATCACATAGCAGTTTCTGAGAATGCTTCTGTCTAGTTTTTATATGAAGATATT
CCATTACTACCATAGGCCTCAAAGCGCTCCAAATATCCACTTGCAGATTCTACAAAAAGAGTGTTTCAA
AACTGCTTTATCAAAAGATAGATTCAACTATGTGAGTTGAATGCACACATAACAAAGTTGTTTCTGAGAA
```

*Figure 2*

- A folder named "figuras", in which Morgan Plotter saves the graphical representation of each analyzed clone
- A .csv file named "tabla_hits". This is the file in which the user specifies the location of each satellite repeat along each one of the sequences to be analyzed and is an input for Morgan Plotter. The file layout must comply with the following format (Fig. 3):

- A column named "hit_id", where an id is provided to each satellite repeat along the clone
- A column named "clone_id", where the name of the clone where each satellite repeat is located
- A column named "identity_percentage", where the similarity between each satellite repeat and the consensus sequence of that specific satellite is displayed
- A column named "hit_start", which specifies the first base of each satellite repeat along the clone
- A column named "hit_end", which specifies the last base of each satellite repeat along the clone

| | A | B | C | D | E | F |
|---|---|---|---|---|---|---|
| 1 | hit_id | clone_id | identity_per | hit_start | hit_end | |
| 2 | 10085H115M | AC007640.13 | 74.251 | 26442 | 26608 | |
| 3 | 100996H130N | AC007640.13 | 74.251 | 29421 | 29255 | |
| 4 | 101130H171N | AC007640.13 | 74.286 | 38861 | 38687 | |
| 5 | 10115H169M | AC007640.13 | 74.405 | 20713 | 20880 | |
| 6 | 101301H75M | AC007640.13 | 74.405 | 33879 | 33712 | |
| 7 | 101471H171N | AC007640.13 | 74.684 | 15371 | 15528 | |
| 8 | 101642H90M | AC007640.13 | 74.843 | 21414 | 21572 | |
| 9 | 101811H171N | AC007640.13 | 74.85 | 25010 | 25176 | |
| 10 | 101982H90M | AC007640.13 | 75 | 27640 | 27807 | |
| 11 | 102084H102N | AC007640.13 | 75 | 41763 | 41600 | |
| 12 | 102226H171N | AC007640.13 | 75 | 47434 | 47604 | |
| 13 | 102397H165N | AC007640.13 | 75 | 48462 | 48632 | |
| 14 | 102568H90M | AC007640.13 | 75 | 50180 | 50349 | |
| 15 | 10263H58M1 | AC007640.13 | 75.301 | 30784 | 30619 | |
| 16 | 1026H171M1 | AC007640.13 | 75.316 | 16236 | 16393 | |
| 17 | 102739H128N | AC007640.13 | 75.449 | 35761 | 35595 | |
| 18 | 10284H85M1 | AC007640.13 | 75.472 | 37638 | 37480 | |
| 19 | 102907H171N | AC007640.13 | 75.595 | 23725 | 23892 | |
| 20 | 103078H85M | AC007640.13 | 75.595 | 31819 | 31652 | |
| 21 | 103248H171N | AC007640.13 | 75.595 | 36103 | 35937 | |
| 22 | 103421H128N | AC007640.13 | 75.595 | 48980 | 49146 | |

*Figure 3*

- A .csv file named "tabla1" which contains the motif to be found, in the 5'-3' orientation, represented as a Position Specific Frequency Matrix. The table should only contain integers, and the *pseudocount* value is 1 [1]. This is an input for Morgan Plotter and Morgan Score (Fig. 4A).
- A .csv file named "tabla2" which contains the motif to be found, in the 3'-5' orientation, represented as a Position Specific Frequency Matrix. The table should only contain integers, and the *pseudocount* value is 1 [1]. This is an input for Morgan Plotter and Morgan Score (Fig. 4B).
- The folder .ipynb_checkpoints, not to be modified by the user.

A

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | |
| 2 | A | 19 | 20 | 8 | 84 | 1 | 96 | 73 | 95 | 54 | 30 | 19 | 24 | |
| 3 | C | 22 | 37 | 16 | 1 | 1 | 1 | 1 | 1 | 1 | 10 | 34 | 34 | |
| 4 | G | 30 | 37 | 8 | 1 | 1 | 1 | 1 | 7 | 21 | 52 | 39 | 34 | |
| 5 | T | 33 | 9 | 73 | 18 | 101 | 6 | 29 | 1 | 28 | 12 | 13 | 12 | |
| 6 | | | | | | | | | | | | | | |

B

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | |
| 2 | A | 12 | 13 | 12 | 28 | 1 | 29 | 6 | 101 | 18 | 73 | 9 | 33 | |
| 3 | C | 34 | 39 | 52 | 21 | 7 | 1 | 1 | 1 | 1 | 8 | 37 | 30 | |
| 4 | G | 34 | 34 | 10 | 1 | 1 | 1 | 1 | 1 | 1 | 16 | 37 | 22 | |
| 5 | T | 24 | 19 | 30 | 54 | 95 | 73 | 96 | 1 | 84 | 8 | 20 | 19 | |
| 6 | | | | | | | | | | | | | | |

*Figure 4*

**MORGAN HUNTER**

Running Morgan Hunter

Open Anaconda Navigator and launch the jupyter notebook (Fig. 5A). A browser menu will be displayed in the default web navigator (Fig. 5B). Make sure that the *Files* tab is selected (upper left corner), open the folder containing the Morgan Hunter Tool and double click Morgan_Hunter.ipynb. The Morgan Hunter interface should appear in a new tab of the web browser (Fig. 5C).

The interface displays a sequence of cells (the first cell is framed by a red rectangle in Fig. 5C) with the commands of the Morgan Hunter algorithm. Cells are run by selecting them (clicking on them) and pressing SHIFT + ENTER. Before a cell is run, brackets to its left appear empty (arrow in Fig. 5C); while the command is running, an asterisk appears inside the brackets; and after the command has been completed, the asterisk is replaced by a number. After a cell is run, the next cell will be automatically selected. If you need to restart running the module, select *Kernel* from the jupyter menu and click on *Restart & Clear Output*.

Before running Morgan Hunter, go to the 10th cell of the algorithm and make sure that the name of the input file ('example_file.txt' in this example, see Fig. 1) is properly written, as shown in the red rectangles in Fig. 6. Go back to the first cell and select it.

Start running the cells of the Morgan Hunter algorithm. The fourth cell requires the user to provide the information of the motif to be found (Fig. 7A, arrow):

- Enter the size of the motif (*Tamaño de la caja*) and press ENTER
- The number of mismatches allowed (*Numero permitido de mismatches*) and press ENTER
- The position of the allowed mismatches separated by commas (*Posiciones para posibles mismatches*) and press ENTER. If no mismatches are allowed, leave the space blank and press ENTER.
- The possible nucleotides of each position along the motif sequence. If two or more nucleotides are allowed in the same position, they should be separated by a pipe (|). Use capital letters only. Press ENTER after each entry (Fig. 7B).

The information provided for the motif search will be displayed on screen (Fig. 7C).
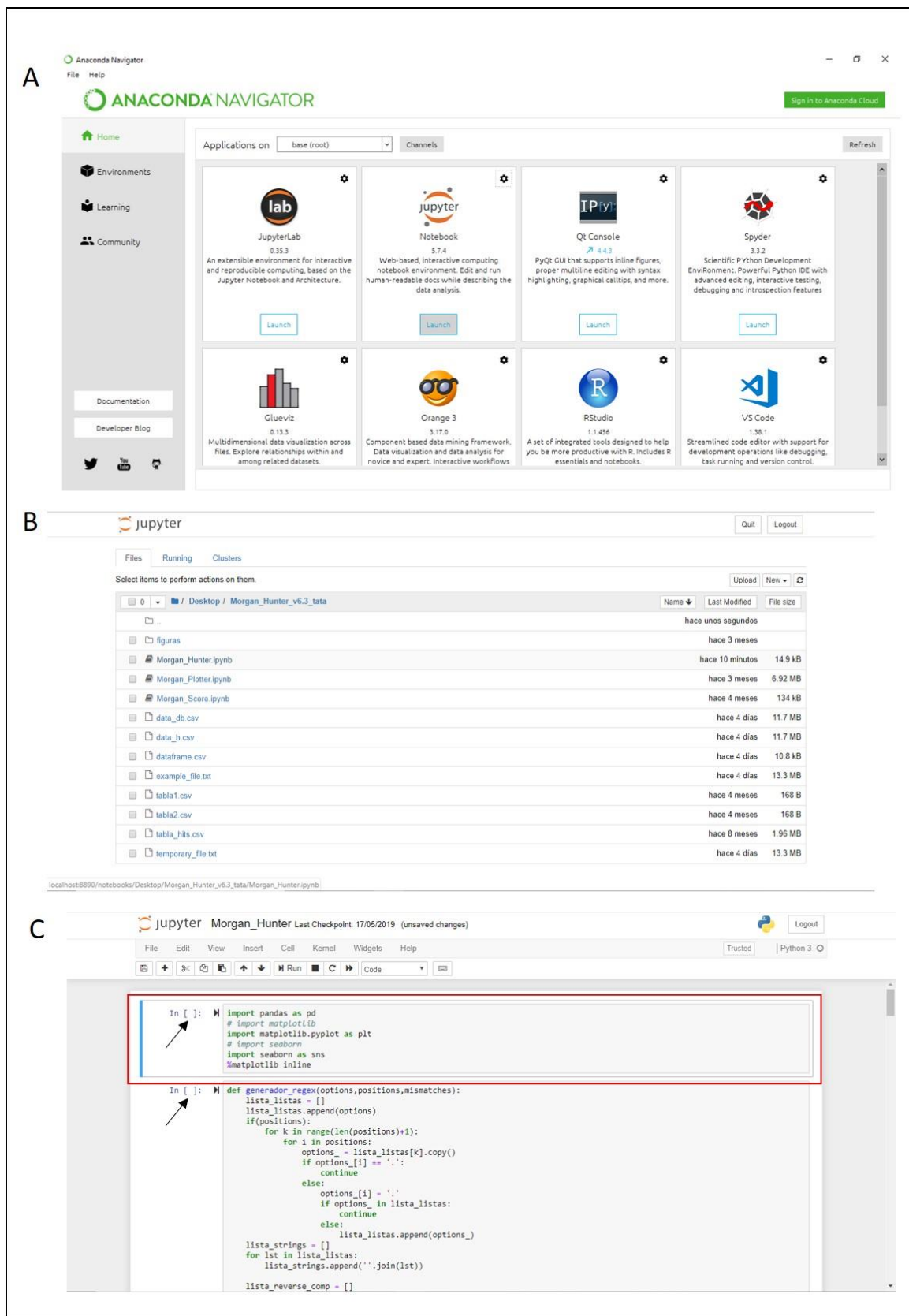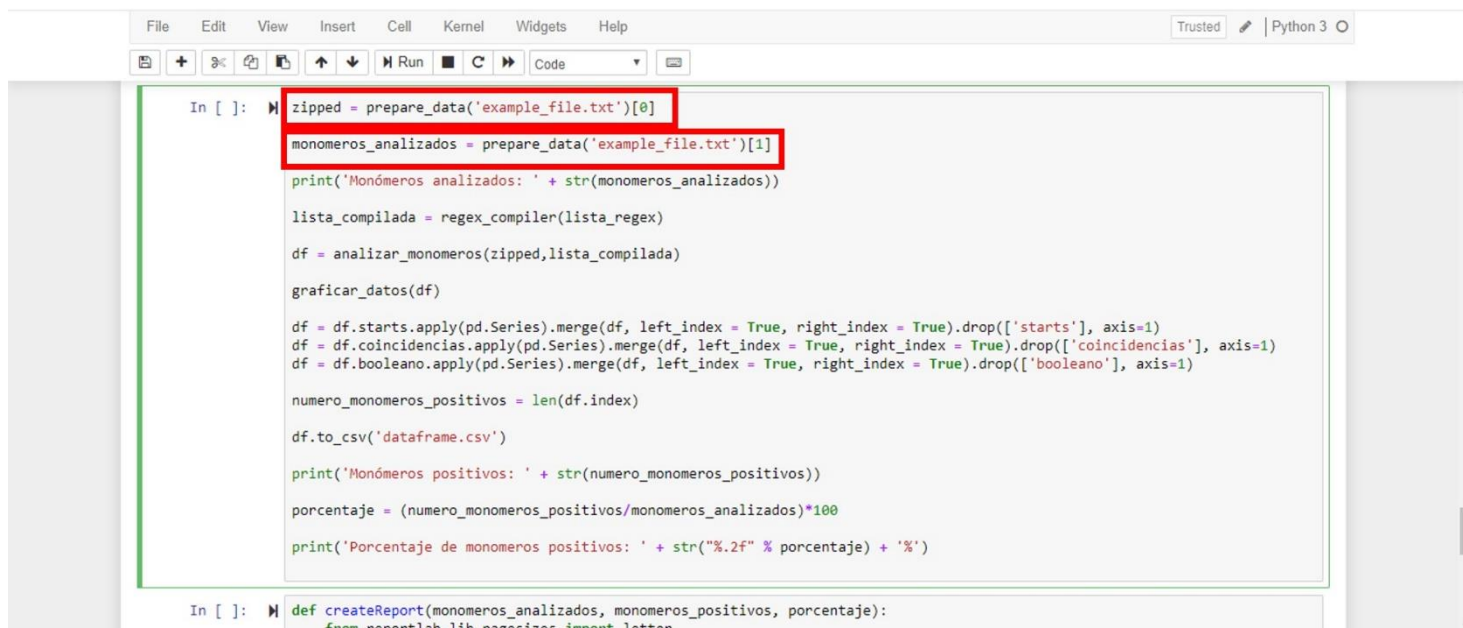
*Figure 5*

*Figure 6*

Morgan Hunter results

Morgan Hunter's outputs are the following:

- The positions (a list of numbers) number (*Monómeros positivos*) and percentage (*Porcentaje de monómeros positivos*) of input sequences that show one or more motifs (displayed on screen below the 10$^{th}$ cell of the algorithm; example shown in Fig. 8)
- A box plot showing the length of the sequences that contain different numbers of motifs (displayed on screen below the 10$^{th}$ cell of the algorithm; example shown in Fig. 8)
- A .csv file named dataframe.csv containing a table showing the coordinates, sequences, and orientation of each motif found along each sequence (saved in the program folder and displayed on screen; example shown in Fig. 9). In this file, each row contains the information of a single sequence:
  - The first column assigns a number to the sequence.
  - The next n columns (where n = the maximum number of motifs found in a single sequence) indicate whether each of the motifs found is in the 5'-3' or 3'-5' orientation (0= 5'-3' and 1=3'-5') (Fig. 9A)
  - The next n columns show the sequence of each motif found (Fig. 9B)
  - The next n columns show the start position of each motif found (Fig. 9C)
  - The last 3 columns contain the sequence ID, the number of motifs found in each sequence, and the size of the entire sequence (Fig. 9C)

  This dataframe is an input file for Morgan Plotter and Morgan Score.
- Two .csv files named data_h.csv and data_db.csv which are required for the following modules of the tool (saved in the program folder; inputs for Morgan Plotter and Morgan Score).

A

```
In [3]:   #Lee archivo con monomeros y crea y devuelve una lista de tuplas de la forma
          #[...,(Nombre,Monomero),...]
          def tuplificar(filename):
              temporary_file =open(filename,'r')
              titulos = []
              monomeros = []
              #Organiza los títulos y monómeros en tuplas
              for line in temporary_file:
                  if line != '\n':
                      if line.find('>') != -1:
                          espacio = line.find(' ')
                          line = line[1:espacio]
                          titulos.append(line)
                      else:
                          line = line.rstrip('\n')
                          monomeros.append(line)

              zipped = list(zip(titulos, monomeros))

              temporary_file.close()
              return zipped
```

```
In [*]:   lista_regex = make_regex_list()
          Tamaño de la caja:  [            ]
```

```
In [ ]:   lista_regex
```

B

```
In [*]:   lista_regex = make_regex_list()

          Tamaño de la caja: 16
          Numero permitido de mismatches: 0
          0
          0123456789101112131415
          Posiciones para posible mismatch(separar con comas):
          (Las opciones van separadas por un pipe '|'.)
          Opciones para la posición 0: A|C|G|T
          Opciones para la posición 1: A|C|G|T
          Opciones para la posición 2: A|C|G|T
          Opciones para la posición 3: G|A
          Opciones para la posición 4: G|A
          Opciones para la posición 5: C
          Opciones para la posición 6: C
          Opciones para la posición 7: A
          Opciones para la posición 8: A
          Opciones para la posición 9: T

          Opciones para la posición 10:  [ G|C         ]
```

```
In [*]:   lista_regex
```

```
In [ ]:   #Compila la lista de expresiones regulares y las regrese como otra lista
          def regex_compiler(lista_regex):
              import re
              lista_compilada = []
              #Compila las expresiones regulares en una lista
              for regex in lista regex:
```

C

```
              temporary_file.close()
              return zipped
```

```
In [4]:   lista_regex = make_regex_list()

          Tamaño de la caja: 16
          Numero permitido de mismatches: 2
          2
          0123456789101112131415
          Posiciones para posible mismatch(separar con comas): 0,15
          (Las opciones van separadas por un pipe '|'.)
          Opciones para la posición 0: A|C
          Opciones para la posición 1: A|C|G|T
          Opciones para la posición 2: A|C|G|T
          Opciones para la posición 3: G|A
          Opciones para la posición 4: G|A
          Opciones para la posición 5: C
          Opciones para la posición 6: C
          Opciones para la posición 7: A
          Opciones para la posición 8: A
          Opciones para la posición 9: T
          Opciones para la posición 10: G|C
          Opciones para la posición 11: G|A
          Opciones para la posición 12: G
          Opciones para la posición 13: A|C|G|T
          Opciones para la posición 14: A|C|G|T
          Opciones para la posición 15: A|C
```

```
In [5]:   lista_regex
```

*Figure 7*

**MORGAN PLOTTER**

Morgan Plotter is opened in the same way as Morgan Hunter. If the Morgan Hunter output files have not been modified, Morgan Plotter can be run by selecting *Cell>Run all* in the jupyter menu or pressing SHIFT + ENTER in each cell, as previously described.

Morgan Plotter results

The numbers that appear under the last cell are the scores of motif/consensus sequence similarity (calculated as in reference [1]) for each motif found (Fig. 10A). These data are also shown in a .csv file named "scores", saved in the program folder. In this file each row represents an analyzed sequence. The first column contains the name of the sequences, and the following contain the score of each of the motifs detected along the corresponding sequence (Fig. 11).

Below the scores, the resulting graphical representations of each analyzed sequence will be displayed on screen, and a .png image of each one of them will be saved in the folder named "figuras". For each analyzed sequence, two horizontal grey bars are displayed, the upper one shows the motifs found in the 5'-3' orientation, while the one below shows the motifs found in the 3'-5' orientation. Vertical lines that cross each grey bar represent the motifs detected by Morgan Hunter. Their score of similarity to the consensus are color coded: dark green = high similarity; dark red = low similarity. The horizontal black lines delimited by vertical black lines (scale bar symbol) within the grey bars represent the satellite repeats provided by the user (Fig. 10B).
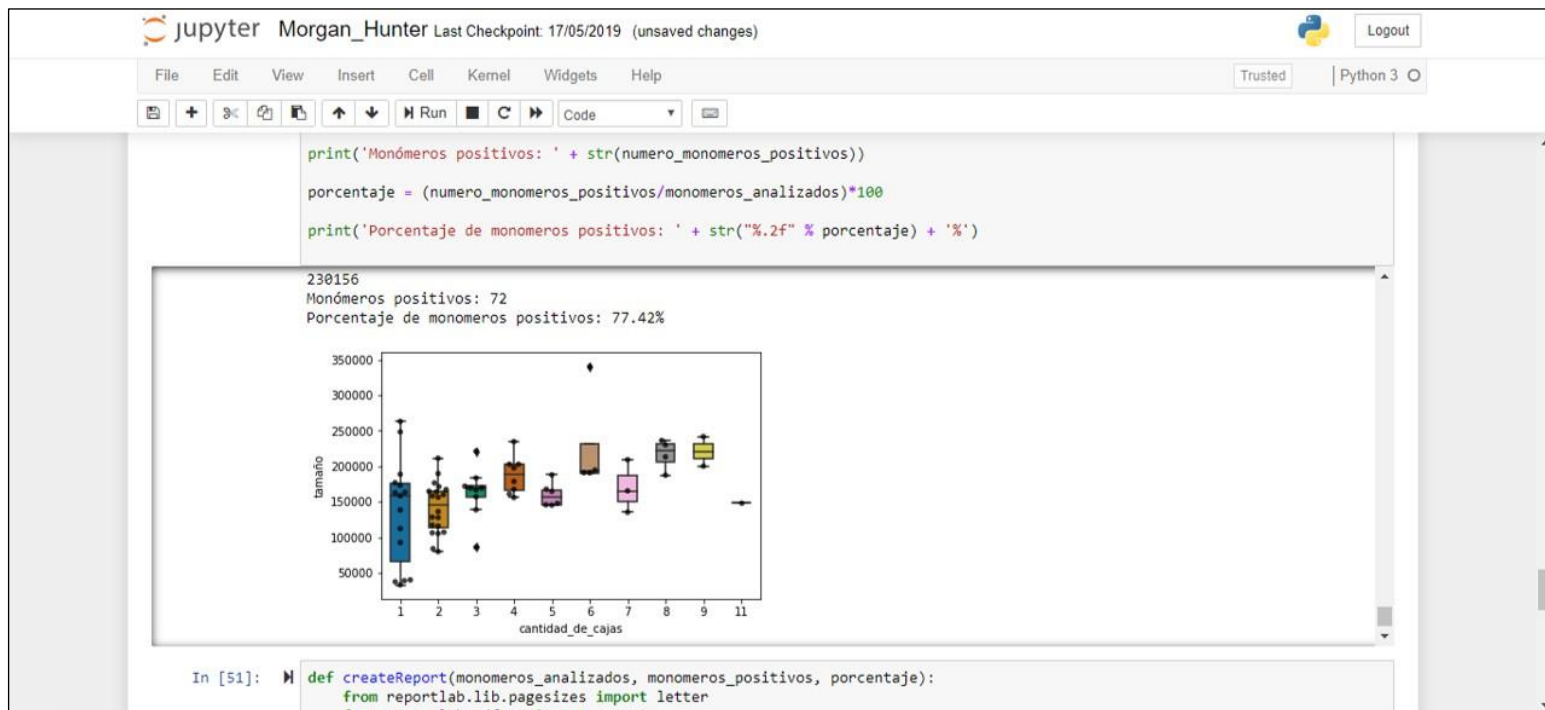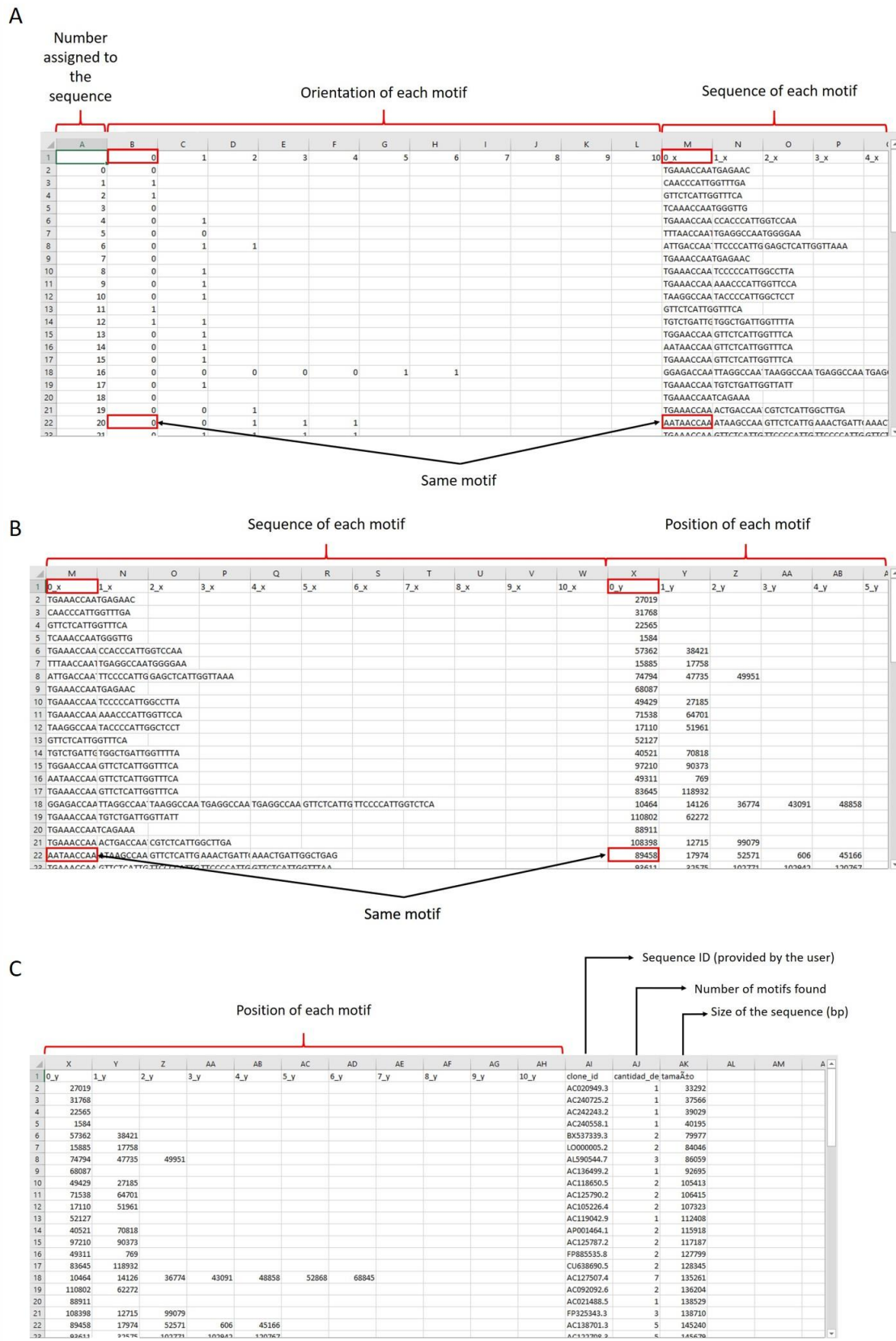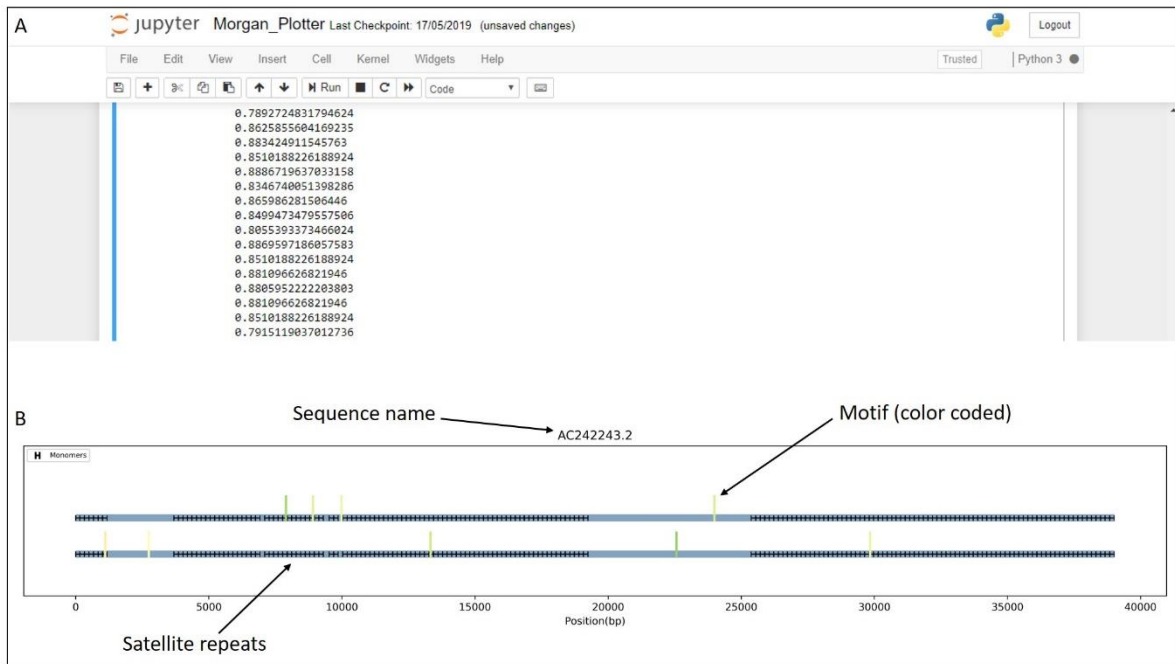


*Figure 8*

*Figure 9*

*Figure 10*

## MORGAN SCORE

Morgan Score is opened and run exactly in the same way as Morgan Plotter. However, it will only generate the consensus similarity scores of each one of the motifs found in the entire input file. This is useful when large amounts of sequences are to be analyzed and no graphical representation is required. With the same input, Morgan Score will run faster than Morgan Plotter. Its output is the same as in Fig. 10A, and it is only displayed on screen.



*Figure 11*

### References

1 Dolfini D, Zambelli F, Pavesi G & Mantovani R (2009) A perspective of promoter architecture from the CCAAT box. *Cell Cycle* 8, 4127–4137.