

Aligning Gene & Protein Sequences

by **Authors:** [Tony Kabilan Okeke](#), [Grace Fan](#)

```
% Your code and output must only contain code and answers for the following
% questions. You must remove or comment out any analysis/code that is not
% for answering these questions.
```

Obtain the gene sequences

Find the hemoglobin subunit alpha-1 genes and proteins for human and Pacific white-sided dolphin. Obtain the sequences of these genes and proteins and download them into the files:

```
hgenefile=bmes.downloadurl('https://www.ncbi.nlm.nih.gov/sviewer/viewer.cgi?tool=portal&save=f
dgenefile=bmes.downloadurl('https://www.ncbi.nlm.nih.gov/sviewer/viewer.cgi?tool=portal&save=f
hptnfile=bmes.downloadurl('https://www.ncbi.nlm.nih.gov/sviewer/viewer.cgi?tool=portal&save=fi
dptnfile=bmes.downloadurl('https://www.ncbi.nlm.nih.gov/sviewer/viewer.cgi?tool=portal&save=fi
```

Read in the fasta files.

```
humgene = fastaread(hgenefile);
dolgene = fastaread(dgenefile);
humptn = fastaread(hptnfile);
dolptn = fastaread(dptnfile);
```

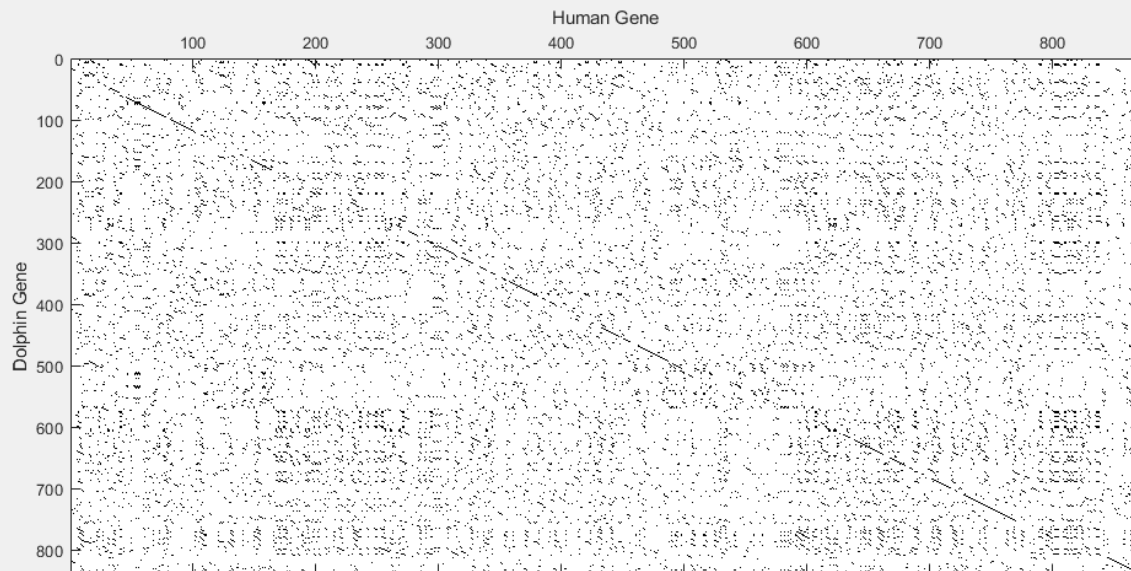
Gene Dot plot

Show sequence dotplot of the the human and dolphin genes. Filter out using a window size of 3, and a match number of 3.

```
seqdotplot(humgene, dolgene, 3, 3)
```

```
Warning: Match matrix has more points than available screen pixels.
Scaling image by factors of 1 in X and 2 in Y.
```

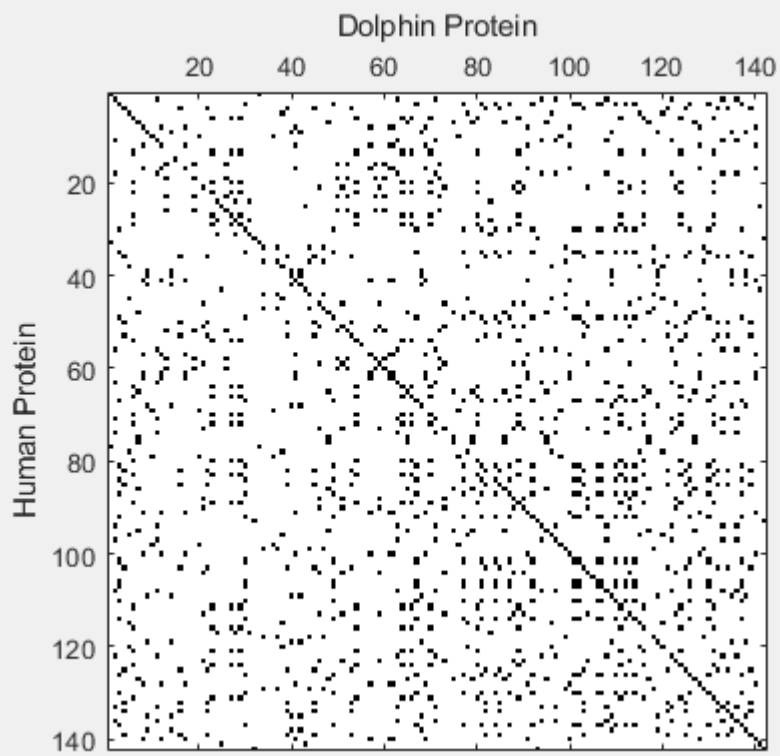
```
xlabel('Human Gene'), ylabel('Dolphin Gene')
```



Protein Dot plot

Show sequence dotplot of the the human and dolphin protein sequences Do not use any filters.

```
seqdotplot(humptn, dolptn)
xlabel('Dolphin Protein'), ylabel('Human Protein')
```



Find out which amino acid has the highest frequency in the **human protein**. Report the single letter amino acid code for it, along with its frequency (as a percentage).

Highest Frequency Amino Acid: A at 14.79%

Perform global alignment of the gene sequences (use default parameters). Show the alignment score and up to first **60** columns of the alignment.

Global Alignment Score: 581.25

```
ans = 3x60 char array
'--CAT---AAACC---C-T---GG--CGCGCTC-G-CGCCCCGGCACTTTCTGGTCC--'
'| | |      ||   | |    ||  ||||| | | | ||||| ||| ||||| ||| |'
'CCCCTCGCGCCCCAGGCATAAAGGCTCGCGCACTGCCAGCCCTGCACGCTTCTGGTCTGT'
```

Local Alignment Score: 614.25

[illegible]

% Percent Identity

Percent Identity: 75.71%

Perform **global alignment** of the protein sequences. Use PAM30 substitution matrix, and an affine gap penalty with gap opening penalty of 8 and gap extension penalty of 1. Show the alignment score and the first **60** columns of the alignment.

Global Alignment Score: 874.00

```
ans = 3x60 char array  
'MVLSPADKTNVKAAGWKVGAHAGEYGAELERMF LSFPTTKTYFPHFDLSHGSAQVKKGHG'  
'||||||| | | :| |: ||||| :||:||| ||| |||:|'|  
'MVLSPADKTNVKGTSKIGNHSAEYGAELERMF INFPS TKTYFSHFDLGHS AQIKKGHG'
```

Percent Identity: 83.80%

Given the DNA sequences ACGTATCGCGTATA and GATGCTCTCGGAAA, find the optimal local alignment and its score. Matches, mismatches, and gaps are scored as +1, 0, and -1, respectively. You can use `swalign()` to solve this problem. You need to specify the 'alphabet' to be nucleotide and provide your own substitution matrix using the 'ScoringMatrix' parameter. See the documentation for `swalign()`.

```
score = 7
align = 3x14 char array
'ACGTATCGCGTATA'
'|:| :||:|:|:|'
'ATG-CTCTCGGAAA'
```