

SeqPlotter: Python package for sequence data analysis and visualization

SeqPlotter is a Python package developed to streamline the visualization and analysis of biological sequence data. It provides seamlessly integration into your Python workflow, with robust data handling capabilities, intuitive visualization, and useful analysis features. From parsing raw sequence files to creating stunning plots and charts.

1. Working with DNA sequences

1.1. Importing DNA module

```
In [1]: # Importing module
from seqplotter.nucl import DNA
```

1.2. Creating DNA objects

```
>>> dna_seq1 = DNA("seqID", "DNA_sequence")
```

Output: retruns a DNA object with seqid and seq pair

```
In [2]: # creating DNA object with seqID and sequence pair
seq1 = DNA("seq1", "GTGTTTGGACTAATAATTGGTCAAGCCTAC")
```

1.3. Sequence length

```
>>> dna_seq1.length()
```

Output: returns sequence length (int) in bp

```
In [3]: # print sequence length
seq1.length()
```

```
Out[3]: 30
```

1.4. Sequence composition

```
>>> dna_seq1.comp()
```

Output: returns a dict object with sequence composition i.e. count of each nucleotides

```
In [4]: # print sequence composition
seq1.comp()
```

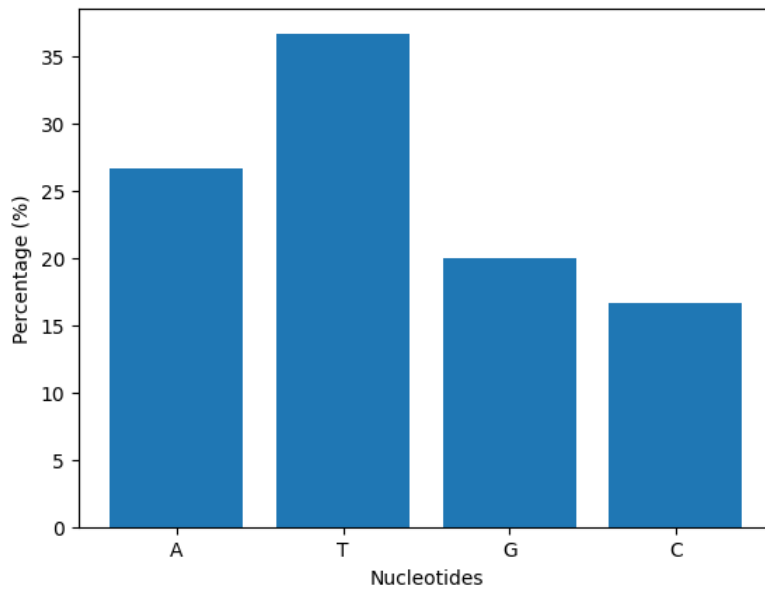
```
Out[4]: {'A': 26.666666666666668,
         'T': 36.666666666666664,
         'G': 20.0,
         'C': 16.666666666666664}
```

1.5 Sequence composition barplot

```
>>> seq1_comp = dna_seq1.comp()
>>> DNA.barplot(seq1_comp)
```

Output: generate barplot showing sequence composition

```
In [5]: # plot sequence composition
comp = seq1.comp()
DNA.barplot(comp)
```



1.6. Sequence slicing

```
>>> # returns nucleotide for the given position
>>> dna_seq1.slice(pos)

>>> # returns nucleotides between start to end position
>>> dna_seq1.slice(start_pos, end_pos)
```

Output: returns a string that contain the sliced sequence

```
In [6]: # slice particular position
seq1.slice(2)
```

```
Out[6]: 'T'
```

```
In [7]: # slice with start and end index
seq1.slice(2,5)
```

```
Out[7]: 'TGTT'
```

1.7. GC content

```
>>> dna_seq1.gc_percent()
```

Output: returns GC percent (string)

```
In [8]: # GC percent
seq1.gc_percent()
```

```
Out[8]: '36.67'
```

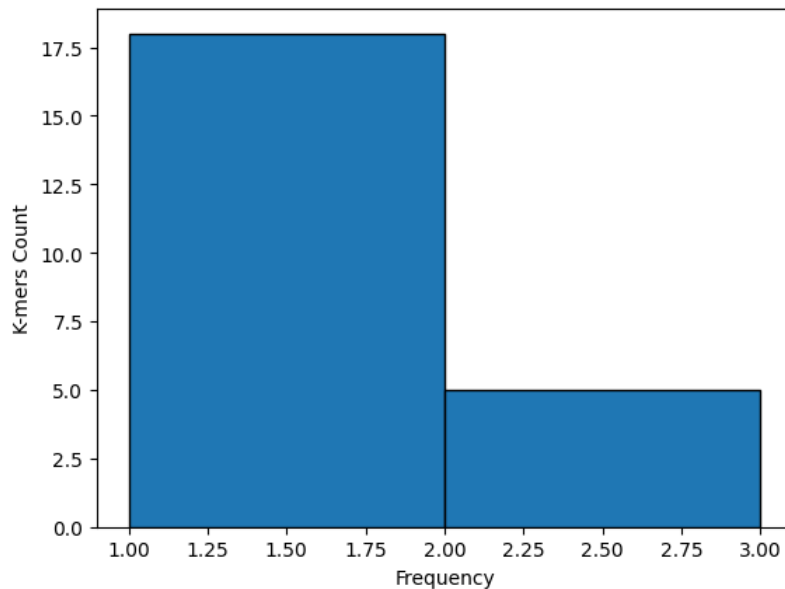
1.8. K-mer abundance plot

```
# generated k-mer plot with default kmer size (k=3)
>>> dna_seq1.kmer_abundance_plot()

# generate k-mer plot with k=6
>>> dna_seq1.kmer_abundance_plot(6)
```

Output: generate k-mer plot (histogram)

```
In [9]: # kmer abundance plot
seq1.kmer_abundance_plot()
```



1.9. Translate gene sequence

```
>>> dna_seq1.translate()
```

Output: translate DNA sequence into amino acid sequence

```
In [10... # translate sequence
seq1.translate()
```

```
Out[10... 'VF*LIIGQAY'
```

1.10. Reading FASTA file (DNA Sequence)

```
>>> fasta_records = DNA.read_fasta("/path/to/fasta/file")
```

Output: returns a list containing DNA object(s)

```
In [11... # reading a FASTA file
data = DNA.read_fasta("seq.fasta")
```

1.11. Using head() and tail() functions

```
>>> # return summary of first 5 sequence (default)
>>> DNA.head()
```

```
>>> # return summary of first 10 sequence
>>> DNA.head(10)
```

```
>>> # return summary of last 5 sequence (default)
>>> DNA.tail()
```

```
>>> # return summary of last 10 sequence
>>> DNA.tail(10)
```

Output: print head and tail summary of parsed FASTA file

```
In [12... # print first five records (by default)
DNA.head()
```

```
seqID    seq
seq1:    [CTTTTCCTTAGTGTTTGGCTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACA]
seq2:    [CTTTTCCTTAGTGTTTGGCTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACA]
seq3:    [CTTTTCCTTAGTGTTTGGCTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACA]
seq4:    [CTTTTCCTTAGTGTTTGGCTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACA]
seq5:    [TCACAGAAGAGCAAGAGGCTCTTGTAGTGAAGTCTTGGAGTGTGATGAAGAAAACTCAGCTGAATTAGG]
```

```
In [13... # print first two records
DNA.head(2)
```

```
seqID    seq
seq1:    [CTTTTCCTTAGTGTTTGGCTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACA]
seq2:    [CTTTTCCTTAGTGTTTGGCTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACA]
```

```
In [14... # print last five records (by default)
DNA.tail()

seqID    seq
seq5:    [TCACAGAAGAGCAAGAGGCTCTTGTAGTGAAGTCTTGAGTGTGTCATGAAGAAAACTCAGCTGAATTAGG]
seq6:    [TCTCAAACCTTTCATCAAGTAAGTAATGATCCCATGATCTCTCTATTTTCTTTTATGTATATAGCA]
seq7:    [TGAGATATGAACACTACTATTTTGAAGTGTAGGATCTTTGAGATTGCACCAACAACGA]
seq8:    [TTCTTGAGAGACTCACCAATTCCTGCTGAGCAAAATCCAAGCTCAAGCCTCACGCAATGTCTGT]
seq9:    [TCATGGTAATAATCAATATCAAATAACATGATTTT]
```

```
In [15... # print last six records
DNA.tail(3)

seqID    seq
seq7:    [TGAGATATGAACACTACTATTTTGAAGTGTAGGATCTTTGAGATTGCACCAACAACGA]
seq8:    [TTCTTGAGAGACTCACCAATTCCTGCTGAGCAAAATCCAAGCTCAAGCCTCACGCAATGTCTGT]
seq9:    [TCATGGTAATAATCAATATCAAATAACATGATTTT]
```

1.12. Count records

```
>>> DNA.count(fasta_records)
```

Output: returns sequence count (int) i.e. number of records

```
In [16... # print record counts
DNA.count(data)
```

```
Out[16... 9
```

1.13. Find record

```
>>> DNA.where(seqid="seq1")
```

Output: returns sequence for a given sequenceID

```
In [17... # find sequence
DNA.where(seqid="seq1")
```

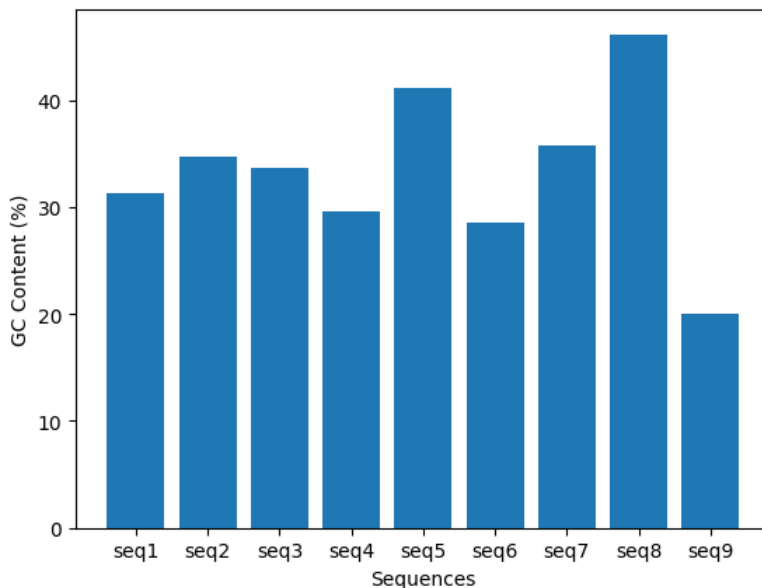
```
Out[17... 'CTTTTCTTTAGTGTGTTTACTAATAATTGGTCAAGCCTACCATTACAACTATGTTCCATTACCAGTACACTTTTCTTTAGTGTGTTTACTAATAATTGGATTACCA
GTACA'
```

1.14. GC percent per sequence plot

```
>>> DNA.gc_plot(fasta_records)
```

Output: generate barplot showing per sequence GC percent

```
In [18... # GC barplot
DNA.gc_plot(data)
```

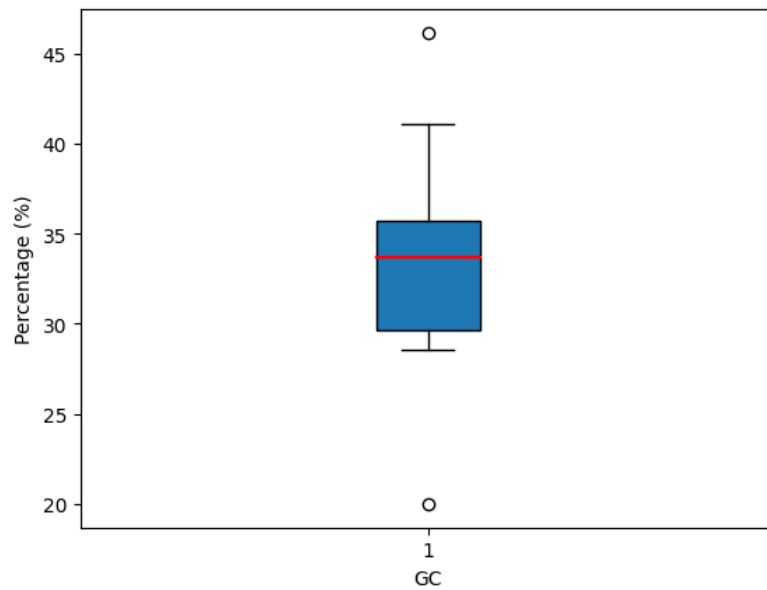


1.15. Boxplot showing GC distribution

```
>>> DNA.gc_distribution_plot(fasta_records)
```

Output: generate boxplot showing overall GC distribution across all the sequences

```
In [19... # GC distribution plot
DNA.gc_distribution_plot(data)
```

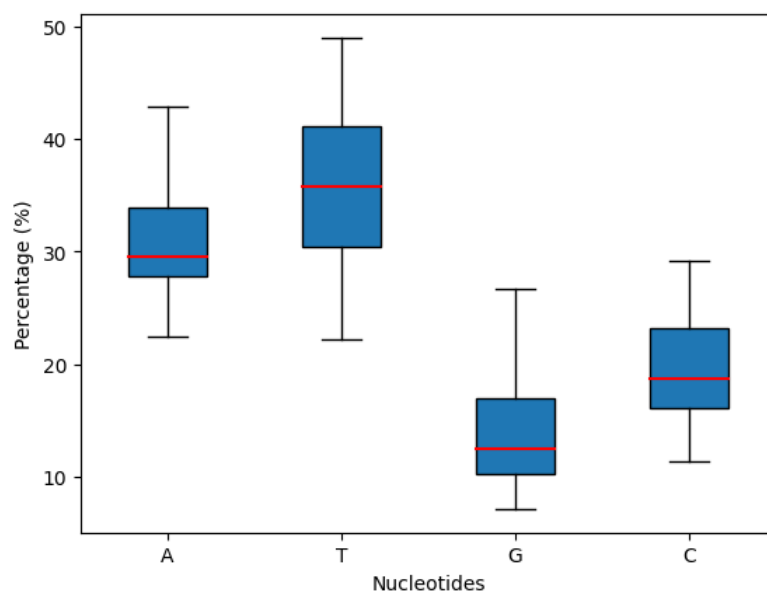


1.16. Boxplot showing nucleotides distribution

```
>>> DNA.nt_distribution_plot(fasta_records)
```

Output: generate boxplot showing overall nucleotides distribution across all the sequences

```
In [20... # Nucleotides distribution plot
DNA.nt_distribution_plot(data)
```

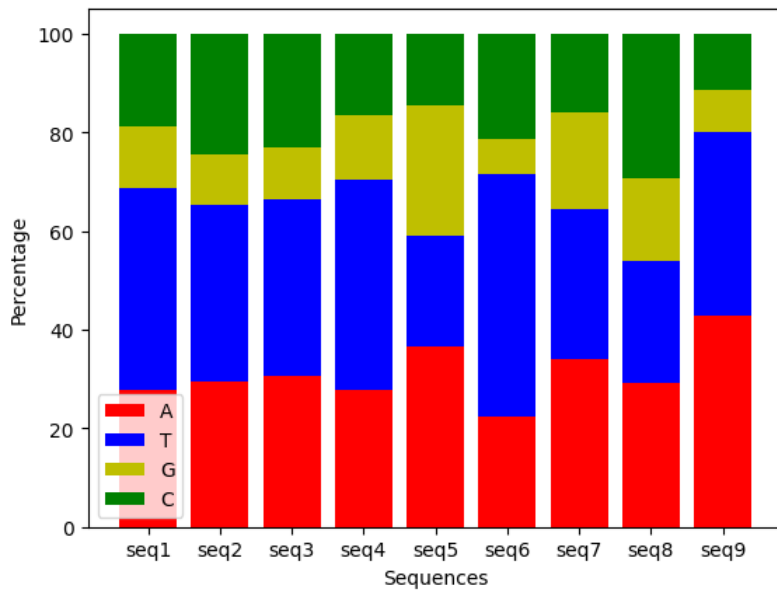


1.17. Stacked barplot showing per sequence base composition

```
>>> DNA.per_sequence_comp(fasta_records)
```

Output: generate stacked barplot showing per sequence nucleotide composition

```
In [21... # Per sequence nucleotides composition
DNA.per_sequence_comp(data)
```

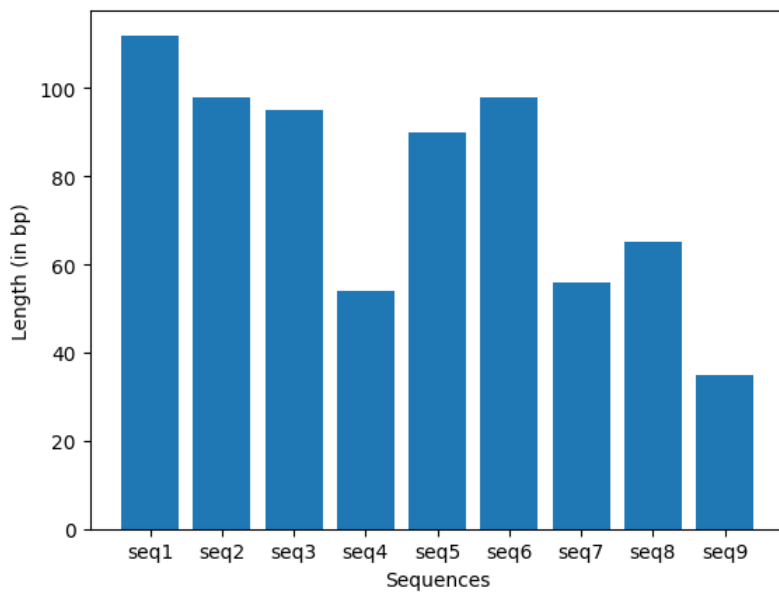


1.18. Barplot showing per sequence length

```
>>> DNA.length_plot(fasta_records)
```

Output: generate barplot showing per sequence length distribution

```
In [22... # Per sequence length
DNA.length_plot(data)
```

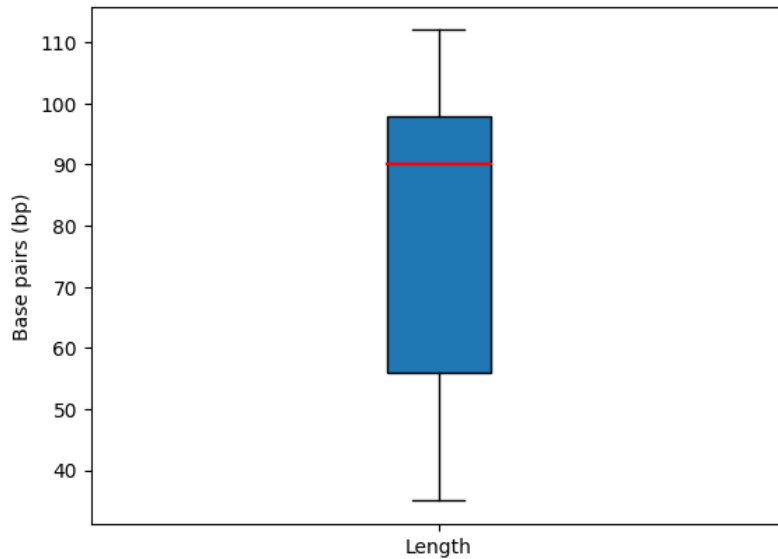


1.19. Boxplot showing overall length distribution

```
>>> DNA.length_distribution_plot(fasta_records)
```

Output: generate boxplot showing overall length distribution across all the sequences

```
In [23... # Boxplot showing length distribution
DNA.length_distribution_plot(data)
```



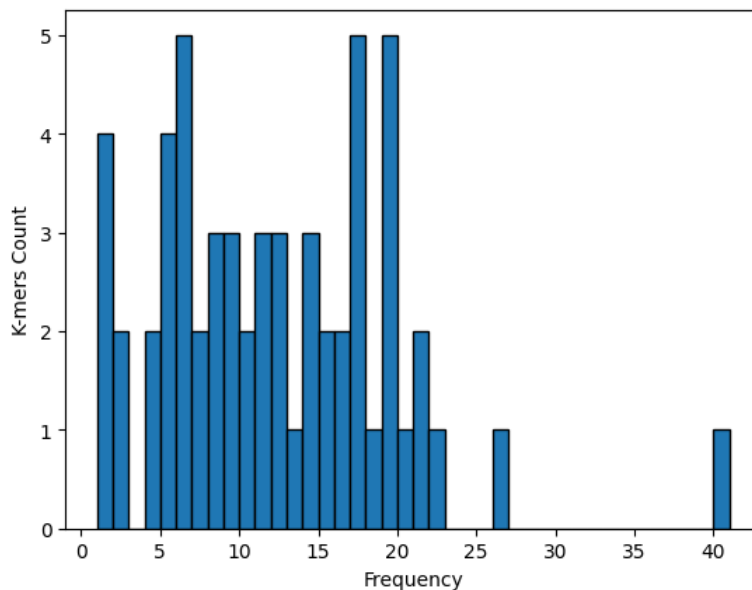
1.20. Overall k-mer abundance pattern

```
# default k-mer size (k=3)
>>> DNA.combined_kmer_abundance_plot(fasta_records)

# kmer size (k=7)
>>> DNA.combined_kmer_abundance_plot(fasta_records, 7)
```

Output: generate combined k-mer abundance plot (histogram) for all the sequence

```
In [24... # Combined kmers distribution
DNA.combined_kmer_abundance_plot(data)
```



1.21. Converting sequence object to feature matrix

```
>>> from seqplotter_analysis import seq2feature
>>> feature_matrix = seq2feature(fasta_records) # with default options (k=3, min_sample=1)

>>> feature_matrix = seq2feature(fasta_records, 6, 2) # with (k=6, min_sample=2)
```

Options:

- k: kmer size
- min_sample: threshold of considering a k-mer as feature

Output: returns a dict object containing following keys ("matrix", "sample", "kmers"), where "kmers" are columns and "sample" are rows

```
< dict{"matrix": numpy_array, "sample": seqID_list, "kmers": kmers_list} >
```

```
In [25... # Importing module
from seqplotter_analysis import seq2feature

# Converting sequence to feature matrix
feature_matrix = seq2feature(data)

# Visualizing feature matrix
import pandas as pd
pd.DataFrame(feature_matrix["matrix"], index=feature_matrix["sample"], columns=feature_matrix["kmers"])
```

```
Out[25...      CTT  TTT  TTC  TCT  TTA  TAG  AGT  GTG  TGT  GTT  ...  GAA  AGA  GAG  GCA  AGG  GCT  CTC  CTG  TGC  ACG
seq1    4   10   3    2    5    2    4    2    3    3  ...    0    0    0    0    0    0    0    0    0    0
seq2    2    5   3    1    4    1    2    1    3    3  ...    0    0    0    0    0    0    0    0    0    0
seq3    2    5   3    1    4    1    3    1    3    3  ...    0    0    0    0    0    0    0    0    0    0
seq4    2    5   2    1    2    1    2    1    2    2  ...    1    0    0    0    0    0    0    0    0    0
seq5    2    0   0    2    2    3    3    2    2    0  ...    7    4    3    1    3    3    3    2    0    0
seq6    3   10   4   12    2    1    2    0    1    0  ...    0    0    0    1    0    0    8    0    0    0
seq7    1    3   0    1    0    1    0    0    1    0  ...    2    2    2    1    1    0    0    1    1    1
seq8    1    0   2    2    0    0    0    0    2    0  ...    0    2    3    2    0    2    3    3    1    1
seq9    0    2   0    0    0    0    0    0    0    0  ...    0    0    0    0    0    0    0    0    0    0
```

9 rows × 54 columns

1.22. Principal component analysis

```
>>> from seqplotter_analysis import PCA
>>> pca_matrix = PCA(feature_matrix) # return pca matrix with first 2 PCs

>>> pca_matrix = PCA(feature_matrix, 3) # return pca matrix with first 3 PCs
```

Output: returns a dict object containing following keys ("matrix", "sample", "components"), where "components" are columns and "sample" are rows

```
< dict{"matrix": numpy_array, "sample": seqID_list, "components": PCA_components} >
```

```
In [26... # Importing module
from seqplotter_analysis import PCA

# Performing PCA
pca_matrix = PCA(feature_matrix)

# Visualizing PCA matrix
import pandas as pd
pd.DataFrame(pca_matrix["matrix"], index=pca_matrix["sample"], columns=pca_matrix["components"])
```

```
Out[26...      PC-1      PC-2
seq1    9.629055    2.243708
seq2    5.947428    3.774627
seq3    6.225862    4.230590
seq4    0.776820    1.919865
seq5   -9.639273    0.927017
seq6    3.503045   -15.066248
seq7   -4.161328    0.626651
seq8   -7.607692    0.437185
seq9   -4.673917    0.906603
```

1.23. Plot scatter plot using PCs

```
>>> from seqplotter_analysis import plot_pca

# Sequence IDs will be treated as class label (default)
>>> plot_pca(pca_matrix)
```

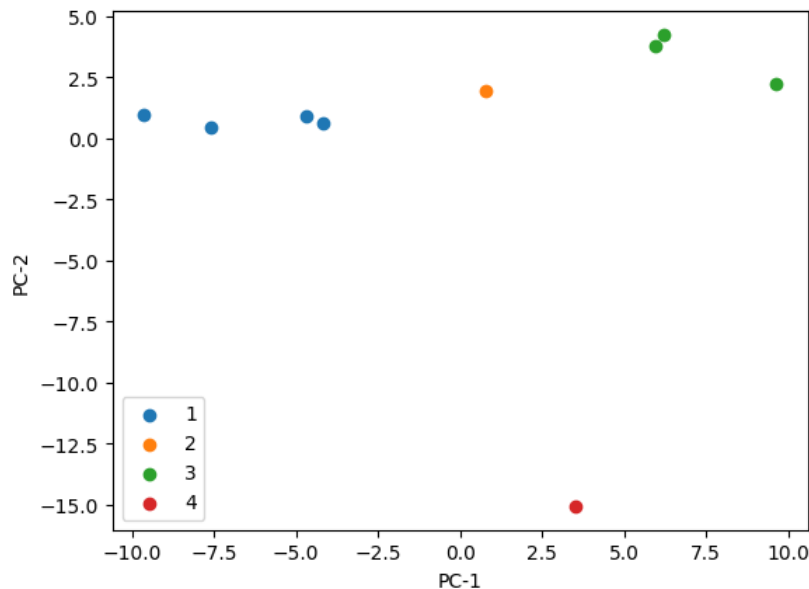


```
# Providing class labels as list
>>> labels = [3, 3, 3, 2, 1, 4, 1, 1, 1]
>>> plot_pca(pca_matrix, labels)
```

Output: generate PCA plot (2D scatter plot using first two components i.e. PC-1 and PC-2)

```
In [27... # Importing module
from seqplotter_analysis import plot_pca

# Plotting 2D PCA plot
labels = [3, 3, 3, 2, 1, 4, 1, 1, 1]
plot_pca(pca_matrix, labels)
```



1.24. Create sequence logo

```
>>> from seqplotter_analysis import seq_logo

# Reading motif sequences file
>>> logo_data = DNA.read_fasta("motif_seq.fasta")

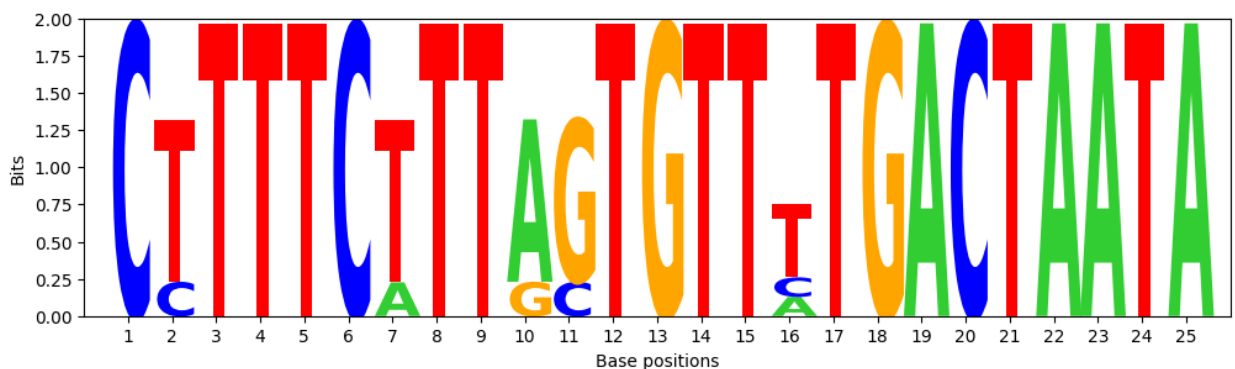
# Generate sequence logo
>>> seq_logo(logo_data, seq_type)
```

Output: generate DNA sequence logo using motif sequences

```
In [28... # Importing module
from seqplotter_analysis import seq_logo

# Reading motif sequences file
data = DNA.read_fasta("dna_motif_seq.fasta")

# Generate sequence logo
seq_logo(data, "DNA")
```



2. Working with Protein sequences

2.1. Importing PROT module

```
In [29... # Importing module
from seqplotter.prot import PROT
```

2.2 Creating PROT object

```
>>> prot_seq1 = PROT("seqID", "protein_sequence")
```

Output: retruns a DNA object with seqid and seq pair

```
In [30... # Creating Protein object with seqID and sequence pair
seq1 = PROT("seq1", "VDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAV")
```

2.3. Sequence length

```
>>> prot_seq1.length()
```

Output: returns sequence length (int) in bp

```
In [31... # print sequence length
seq1.length()
```

```
Out[31... 35
```

2.4. Sequence composition

```
>>> prot_seq1.comp()
```

Output: returns a dict object with sequence composition i.e. count of each amino acids

```
In [32... # print sequence composition
seq1.comp()
```

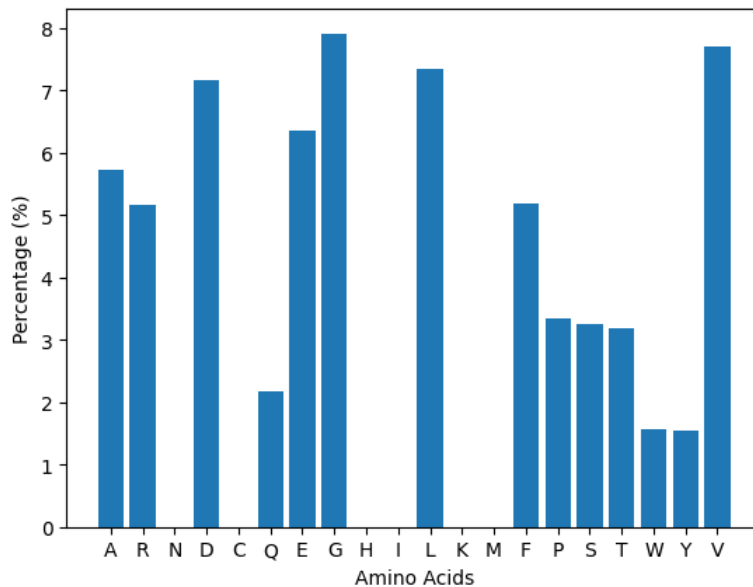
```
Out[32... {'A': 5.714285714285714,
'R': 5.166051660516605,
'N': 0.0,
'D': 7.163266076755573,
'C': 0.0,
'Q': 2.1718543403148107,
'E': 6.353851345091472,
'G': 7.909936022986237,
'H': 0.0,
'I': 0.0,
'L': 7.342245635482005,
'K': 0.0,
'M': 0.0,
'F': 5.188382310331203,
'P': 3.3327849177003355,
'S': 3.260373880823663,
'T': 3.1947334484658114,
'W': 1.567453020447904,
'Y': 1.5536341292308162,
'V': 7.701922950367426}
```

2.5. Sequence composition plot

```
>>> prot_comp = prot_seq1.comp()
>>> PROT.barplot(prot_comp)
```

Output: generate barplot showing sequence composition

```
In [33... # plot sequence composition
comp = seq1.comp()
PROT.barplot(comp)
```



2.6. Molecular weight in KDa

```
>>> prot_seq1.mol_weight()
```

Output: return molecular weight (int) of a protein sequence in KDa

```
In [34... # molecular weight
seq1.mol_weight()
```

```
Out[34... 3.87
```

2.7. Sequence slicing

```
>>> # returns amino acid for the given position
>>> prot_seq1.slice(pos)

>>> # returns amino acids between start to end position
>>> prot_seq1.slice(start_pos, end_pos)
```

Output: returns a string that contain the sliced sequence

```
In [35... # slice particular position
seq1.slice(6)
```

```
Out[35... 'G'
```

```
In [36... # slice with start and end index
seq1.slice(4,8)
```

```
Out[36... 'VGGEA'
```

2.8. Reading FASTA file (Protein sequence)

```
>>> fasta_records = PROT.read_fasta("/path/to/fasta/file")
```

Output: returns a list containing PROT object(s)

```
In [37... # reading a FASTA file
data = PROT.read_fasta("prot_seq.fasta")
```

2.9. Using head() and tail() functions

```
>>> # return summary of first 5 sequence (default)
>>> PROT.head()

>>> # return summary of first 10 sequence
>>> PROT.head(10)

>>> # return summary of last 5 sequence (default)
>>> PROT.tail()
```

```
>>> # return summary of last 10 sequence
>>> PROT.tail(10)
```

Output: print head and tail summary of parsed FASTA file

```
In [38... # Head
PROT.head(5)

seqID  seq
seq1:  [MVHLTPEEKSAVTALWGKVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVKAHGKKVLG]
seq2:  [MLPFWKRLLYAAVIAGALVGADAQFWKTAGTAGSIQDSVKHYNRNEPKFPIDDSYDIVDSAGVARGDLPP]
seq3:  [TGQKGDRGEPGLNGLPGNPGQKGEPGRAGATGKPGLLGPPGPPGGGRGTPGPPGPKGPRGYVGAPGPQGL]
seq4:  [LPGATGEPGKPALCDLSLIEPLKGDGKYPGAPGAKGVQGFKAEGLPGIPGPKGEFGFKGEKGLSGAPGN]
```

```
In [39... # Tail
PROT.tail()

seqID  seq
seq1:  [MVHLTPEEKSAVTALWGKVNDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAMGNPKVKAHGKKVLG]
seq2:  [MLPFWKRLLYAAVIAGALVGADAQFWKTAGTAGSIQDSVKHYNRNEPKFPIDDSYDIVDSAGVARGDLPP]
seq3:  [TGQKGDRGEPGLNGLPGNPGQKGEPGRAGATGKPGLLGPPGPPGGGRGTPGPPGPKGPRGYVGAPGPQGL]
seq4:  [LPGATGEPGKPALCDLSLIEPLKGDGKYPGAPGAKGVQGFKAEGLPGIPGPKGEFGFKGEKGLSGAPGN]
```

2.10. Count records

```
>>> PROT.count(fasta_records)
```

Output: returns sequence count (int) i.e. number of records

```
In [40... # Sequence count
PROT.count(data)
```

```
Out[40... 4
```

2.11. Find record

```
>>> PROT.where(seqid="seq2")
```

Output: returns sequence for the given sequenceID

```
In [41... # Find sequence
PROT.where(seqid="seq2")
```

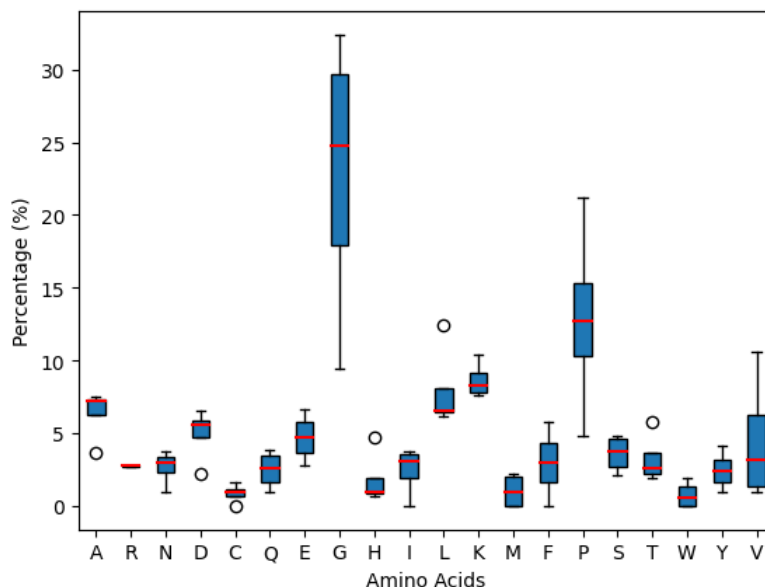
```
Out[41... 'MLPFWKRLLYAAVIAGALVGADAQFWKTAGTAGSIQDSVKHYNRNEPKFPIDDSYDIVDSAGVARGDLPPKNCTAGYAGCVPKCIAEKGNRGLPGPLGPTGLKGEMG
FPGMEGPSGDKGQKGDPGYPGQRGDKGERGSPGLHGQAGVPGVQGPAGNPGAPGINGKDGCDGQDGIPLGLEGLSGMPGPRGYAGQLGSKGEKGEPAKENGDYA'
```

2.12. Amino acids distribution plot

```
>>> PROT.aa_distribution_plot(fasta_records)
```

Output: generate boxplot showing overall amino acids distribution across all the sequences

```
In [42... # Boxplot aa distribution plot
PROT.aa_distribution_plot(data)
```

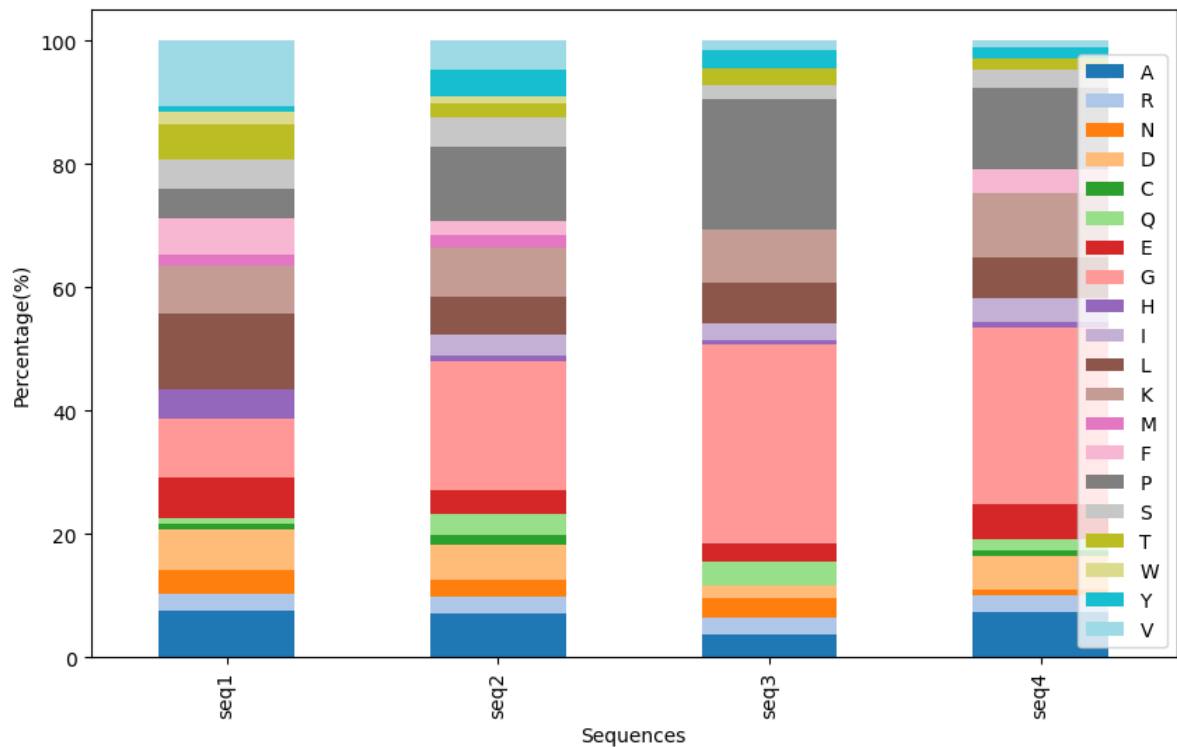


2.13. Stacked barplot showing per sequence amino acid composition

```
>>> PROT.per_sequence_comp(fasta_records)
```

Output: generate stacked barplot showing per sequence amino acids composition

```
In [43... # Per sequence nucleotides composition
PROT.per_sequence_comp(data)
```

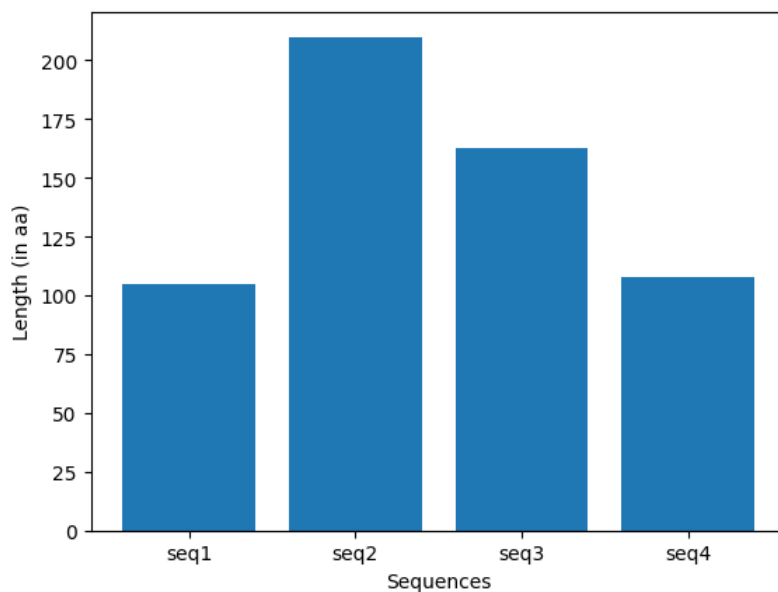


2.14. Barplot showing per sequence length

```
PROT.length_plot(fasta_records)
```

Output: generate barplot showing per sequence length distribution

```
In [44... # Per sequence length
PROT.length_plot(data)
```

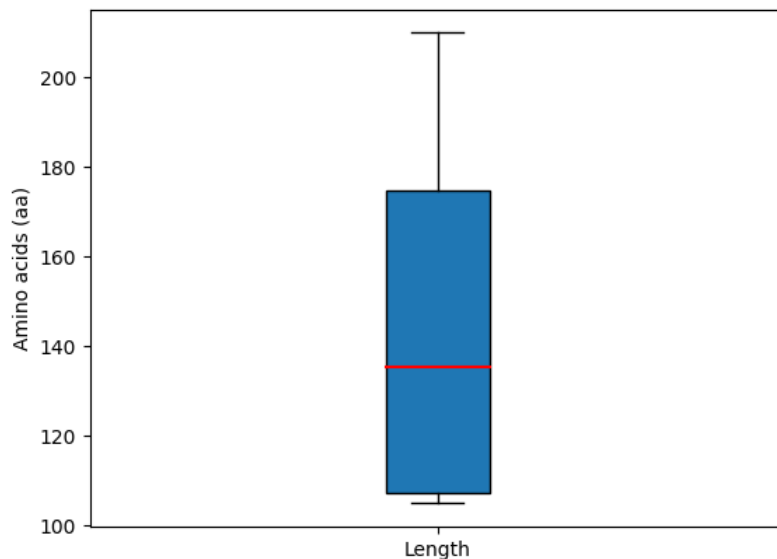


2.15. Boxplot showing overall length distribution

```
>>> PROT.length_distribution_plot(fasta_records)
```

Output: generate boxplot showing overall length distribution across all the sequences

```
In [45... PR0T.length_distribution_plot(data)
```



2.16. Converting sequence object to feature matrix

```
>>> from seqplotter_analysis import seq2feature
>>> feature_matrix = seq2feature(fasta_records) # with default options (k=3, min_sample=1)

>>> feature_matrix = seq2feature(fasta_records, 6, 2) # with (k=6, min_sample=2)
```

Options:

- k: kmer size
- min_sample: threshold of considering a k-mer as feature

Output: returns a dict object containing following keys ("matrix", "sample", "kmers"), where "kmers" are columns and "sample" are rows

```
< dict{"matrix": numpy_array, "sample": seqID_list, "kmers": kmers_list} >
```

```
In [46... # Importing module
from seqplotter_analysis import seq2feature

# Converting sequence to feature matrix
feature_matrix = seq2feature(data)

# Visualizing feature matrix
import pandas as pd
pd.DataFrame(feature_matrix["matrix"], index=feature_matrix["sample"], columns=feature_matrix["kmers"])
```

```
Out[46...      RLL  GDL  DLS  GNP  GKK  DGL  LKG  KGT  FWK  AGA  ...  DAG  PGY  GYP  YPG  GAK  AKG  GFK  FKG  GRD  RDG
seq1    1    1    1    1    1    1    1    1    0    0  ...    0    0    0    0    0    0    0    0    0
seq2    1    1    0    1    0    0    1    0    2    1  ...    0    0    0    0    0    0    0    0    0
seq3    0    0    0    1    1    1    1    1    0    1  ...    2    2    1    1    0    0    0    0    0
seq4    0    0    1    0    0    0    1    0    0    0  ...    0    0    2    2    2    2    2    2    2
```

4 rows × 84 columns

2.17 Principal component analysis

```
>>> from seqplotter_analysis import PCA
>>> pca_matrix = PCA(feature_matrix) # return pca matrix with first 2 PCs

>>> pca_matrix = PCA(feature_matrix, 3) # return pca matrix with first 3 PCs
```

Output: returns a dict object containing following keys ("matrix", "sample", "components"), where "components" are columns and "sample" are rows

```
< dict{"matrix": numpy_array, "sample": seqID_list, "components": PCA_components} >
```

```
In [47... # Importing module
from seqplotter_analysis import PCA

# Performing PCA
pca_matrix = PCA(feature_matrix)

# Visualizing PCA matrix
import pandas as pd
pd.DataFrame(pca_matrix["matrix"], index=pca_matrix["sample"], columns=pca_matrix["components"])
```

```
Out[47...
      PC-1    PC-2
seq1  5.069106 -3.616282
seq2  2.496495 -2.811089
seq3 -11.398262 -0.046961
seq4  3.832661  6.474332
```

2.18. Plot scatter plot using PCs

```
>>> from seqplotter_analysis import plot_pca

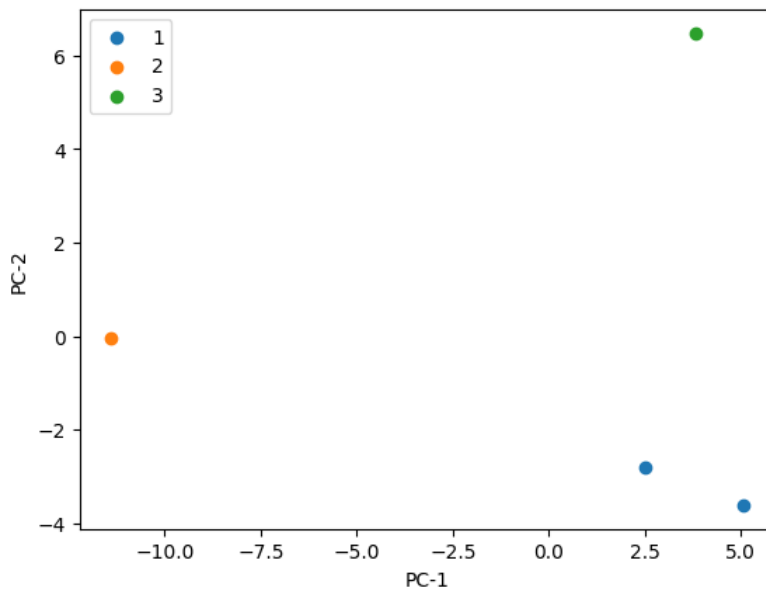
# Sequence IDs will be treated as class label (default)
>>> plot_pca(pca_matrix)

# Providing class labels as list
>>> labels = [3, 3, 3, 2, 1, 4, 1, 1, 1]
>>> plot_pca(pca_matrix, labels)
```

Output: generate PCA plot (2D scatter plot using first two components i.e. PC-1 and PC-2)

```
In [48... # Importing module
from seqplotter_analysis import plot_pca

# Plotting 2D PCA plot
labels = [1, 1, 2, 3]
plot_pca(pca_matrix, labels)
```



2.19. Create sequence logo

```
>>> from seqplotter_analysis import seq_logo

# Reading motif sequences file
>>> logo_data = PROT.read_fasta("motif_seq.fasta")

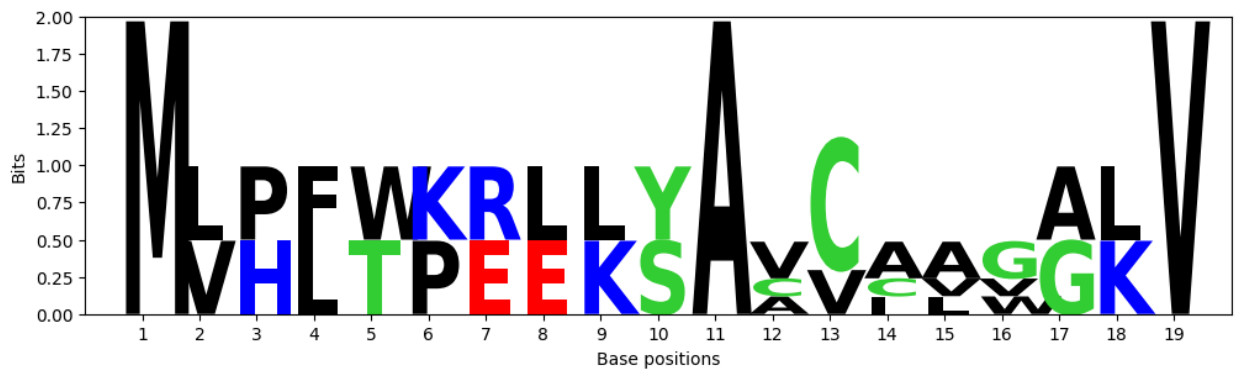
# Generate sequence logo
>>> seq_logo(logo_data, seq_type)
```

Output: generate protein sequence logo using motif sequences

```
In [49... # Importing module
from seqplotter_analysis import seq_logo
```

```
# Reading motif sequences file
data = PROT.read_fasta("prot_motif_seq.fasta")

# Generate sequence logo
seq_logo(data, "PROT")
```



2.20. Generate hydrophobicity plot

```
>>> from seqplotter_analysis import hydrophobicity_plot

# Amino acid sequence for seq2 seqID
>>> seq = PROT.where(seqid="seq2")

# Create hydrophobicity plot
>>> hydrophobicity_plot(seq, window_size)
```

Output: generate hydrophobicity plot for a given protein sequence

```
In [50... from seqplotter_analysis import hydrophobicity_plot

# Amino acid sequence for seq2 seqID
seq = PROT.where(seqid="seq2")

# Create hydrophobicity plot
hydrophobicity_plot(seq, 9)
```

