

新人課題

津嶋 佑旗

2019 年 4 月 19 日

1 配列情報処理

1.1 ヒト 21 番染色体の先頭コンティグの DNA 配列を秋山研サーバから取得し、A/T/G/C の各塩基数を出力せよ。

以下のように実行する。実行環境は Python2 で、ghostgw 上での動作を確認している。

```
1 python 1_1.py /mnt/fs/ohue/newcomer/NT_113952.1.fasta
```

ソースコード 1 実行方法

プログラム本体は以下の通りである。fasta ファイルの塩基配列について、行ごとに Python の機能で各文字の出現回数をカウントする。

```
1 #2019/04/11
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7
8 #file open
9 args = sys.argv
10 if len(args) != 2:
11     sys.stderr.write('Usage: python 1_1.py ./NT_113952.1.fasta\n')
12     sys.exit()
13 path = args[1]
14 try:
15     file = open(path)
16 except:
17     sys.stderr.write('Error in opening file.\n')
18     sys.exit()
19 lines = file.readlines()
20 length = len(lines)
21
22 #counter(num of A/T/G/C)
23 noA = 0
24 noT = 0
25 noG = 0
26 noC = 0
27 #Process with each line
28 #Skip Header
29 for l in range(1, length):
30     noA = noA + lines[l].count('A')
31     noT = noT + lines[l].count('T')
32     noG = noG + lines[l].count('G')
33     noC = noC + lines[l].count('C')
34 print("A:"+str(noA))
35 print("T:"+str(noT))
36 print("G:"+str(noG))
37 print("C:"+str(noC))
38 file.close()
```

ソースコード 2 1_1.py

実行した結果が以下の通りである.

```
1 A:59200
2 T:56195
3 G:34714
4 C:34246
```

ソースコード 3 実行結果

1.2 NT_113952.1.fasta の逆相補鎖を生成するプログラムを作成せよ.

以下のように実行する. 実行環境は Python2 で, ghostgw 上での動作を確認している.

```
1 python 1_2.py /mnt/fs/ohue/newcomer/NT_113952.1.fasta
```

ソースコード 4 実行方法

プログラム本体は以下の通りである.

```
1 #2019/04/05
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7 from rev import rev
8 from combine import combine
9
10 #file open
11 args = sys.argv
12 if len(args) != 2:
13     sys.stderr.write('Usage: python 1_2.py ./NT_113952.1.fasta\n')
14     sys.exit()
15 path = args[1]
16 try:
17     file = open(path)
18 except:
19     sys.stderr.write('Error in opening file.\n')
20     sys.exit()
21 lines = file.readlines()
22
23 #Skip Header
24 lines.pop(0)
25 text = combine(lines)
26 #Reverse
27 reverse = rev(text)
28 print(reverse)
29
30 file.close()
```

ソースコード 5 1_2.py

```
1 #2019/04/05
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4
5 def rev(text):
6     ans = ''
7     for char in text:
8         if char == 'A':
9             ans = 'T' + ans
10        elif char == 'T':
11            ans = 'A' + ans
12        elif char == 'G':
13            ans = 'C' + ans
14        elif char == 'C':
15            ans = 'G' + ans
16    return ans
```

ソースコード 6 rev.py

```

1 #2019/04/05
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4
5 def combine(file):
6     ans = ''
7     for line in file:
8         ans = ans + line.strip()
9     return ans

```

ソースコード 7 combine.py

実行結果は result_1.2.txt として別添する.

1.3 ウィンドウ幅 w , ステップ幅 s で, ウィンドウ内の GC 含量を出力するプログラムを作成し, NT_113952.1.fasta に $w = 1000$, $s = 300$ で適用した結果を gnuplot 等 Excel 以外のツールでプロットせよ.

以下のように実行する. 実行環境は Python2 で, ghostgw 上での動作を確認している.

```

1 python 1_3.py /mnt/fs/ohue/newcomer/NT_113952.1.fasta 1000 300 > data.txt

```

ソースコード 8 実行方法

プログラム本体は以下の通りである. また, 7 も使用する.

```

1 #2019/04/11
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7 from combine import combine
8
9 #file open
10 args = sys.argv
11 if len(args) != 4:
12     sys.stderr.write('Usage: python 1_3.py ./NT_113952.1.fasta 1000 300\n')
13     sys.exit()
14 path = args[1]
15 try:
16     file = open(path)
17 except:
18     sys.stderr.write('Error in opening file.\n')
19     sys.exit()
20 lines = file.readlines()
21
22 #Skip Header
23 lines.pop(0)
24 text = combine(lines)
25 length = len(text)
26
27 try:
28     w = int(args[2])
29     s = int(args[3])
30 except ValueError as e:
31     sys.stderr.write('Incorrect w or s in command line args\n')
32     sys.exit()
33
34 GC = []
35 bottom = 0
36 while bottom < length:
37     part = text[bottom:bottom+w]
38     if len(part) == 0:
39         break
40     GC.append((part.count('C') + part.count('G'))/float(len(part)))
41     bottom = bottom + s
42 #Export for gnuplot
43 for i in range(len(GC)):
44     print(str(i*300)+' '+str(GC[i]))
45
46 file.close()

```

実行結果を gnuplot で描画した結果以下のようになった。

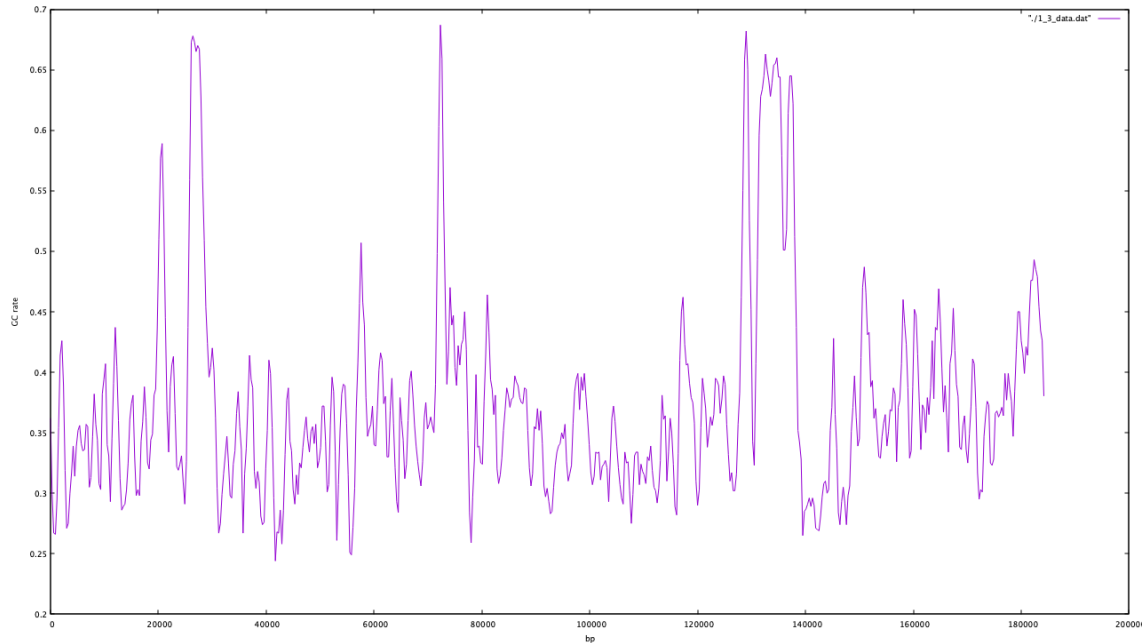


図1 グラフ

1.4 引数で与えられた部分配列を検索し、何文字目に現れたかを表示するプログラムを作成せよ (逆相補鎖上も検索すること). 部分配列を GAATTC および ATG とし、NT_113952.1.fasta に適用せよ.

以下のように実行する. 実行環境は Python2 で, ghostgw 上での動作を確認している.

```
1 python 1_4.py /mnt/fs/ohue/newcomer/NT_113952.1.fasta GAATTC
2 python 1_4.py /mnt/fs/ohue/newcomer/NT_113952.1.fasta ATG
```

ソースコード 10 実行方法

プログラム本体は以下の通りである. また, 6 と 7 も使用する.

```
1 #2019/04/11
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7 from combine import combine
8 from rev import rev
9 from findall import findall
10
11 #file open
12 args = sys.argv
13 if len(args) != 3:
14     sys.stderr.write('Usage: python 1_4.py ./NT_113952.1.fasta GAATTC\n')
15     sys.exit()
16 path = args[1]
17 try:
18     file = open(path)
```

```

19 except:
20     sys.stderr.write('Error in opening file.\n')
21     sys.exit()
22 lines = file.readlines()
23
24 lines.pop(0)
25 text = combine(lines)
26 reverse = rev(text)
27 index1 = findall(text, args[2])
28 index2 = findall(reverse, args[2])
29 print(args[2]+' in STRING is below:')
30 print(index1)
31 print(args[2]+' in REVERSE STRING is below:')
32 print(index2)
33 file.close()

```

ソースコード 11 1_4.py

```

1 #2019/04/08
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4
5 def findall(text, query):
6     index = []
7     data = -1
8     while True:
9         data = text.find(query, data+1)
10        if data == -1:
11            break
12        index.append(data)
13    return index

```

ソースコード 12 findall.py

実行した結果が以下の通りである。ただし、ATG の結果については極端に長いため result_1_4.ATG.txt として別添する。

```

1 GAATTC in STRING is below:
2 [395, 4319, 5056, 5334, 10567, 14092, 14296, 15029, 19284, 22535, 25293, 25897, 30824, 34194,
3 34625, 37333, 56600, 57464, 58202, 60553, 65241, 66725, 79939, 80725, 82012, 86322, 92914,
4 94554, 96877, 97628, 99894, 102214, 103133, 120291, 121414, 121612, 124893, 124901, 125389,
5 151452, 151874, 156137, 158738, 166310, 175826, 177983, 180004, 182341]
6
7 GAATTC in REVERSE STRING is below:
8 [2008, 4345, 6366, 8523, 18039, 25611, 28212, 32475, 32897, 58960, 59448, 59456, 62737, 62935,
9 64058, 81216, 82135, 84455, 86721, 87472, 89795, 91435, 98027, 102337, 103624, 104410,
10 117624, 119108, 123796, 126147, 126885, 127749, 147016, 149724, 150155, 153525, 158452,
11 159056, 161814, 165065, 169320, 170053, 170257, 173782, 179015, 179293, 180030, 183954]

```

ソースコード 13 実行結果 (GAATTC)

1.5 NT_113952.1.fasta を 6 つの読み枠でアミノ酸配列に変換せよ。Stop コドンはアンダースコアで表示すること。

以下のように実行する。実行環境は Python2 で、ghostgw 上での動作を確認している。

```

1 python 1_5.py /mnt/fs/ohue/newcomer/NT_113952.1.fasta

```

ソースコード 14 実行方法

プログラム本体は以下の通りである。また、6 と 7 も使用する。

```

1 #2019/04/11
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7 from combine import combine
8 from rev import rev
9 from decode import decode

```

```

10
11 #file open
12 args = sys.argv
13 if len(args) != 2:
14     sys.stderr.write('Usage: python 1_5.py ./NT_113952.1.fasta\n')
15     sys.exit()
16 path = args[1]
17 try:
18     file = open(path)
19 except:
20     sys.stderr.write('Error in opening file.\n')
21     sys.exit()
22 lines = file.readlines()
23
24 lines.pop(0)
25 text = combine(lines)
26 reverse = rev(text)
27 print('PEPTIDE 1')
28 print(decode(text,0))
29 print('PEPTIDE 2')
30 print(decode(text,1))
31 print('PEPTIDE 3')
32 print(decode(text,2))
33 print('PEPTIDE 4')
34 print(decode(reverse,0))
35 print('PEPTIDE 5')
36 print(decode(reverse,1))
37 print('PEPTIDE 6')
38 print(decode(reverse,2))
39 file.close()

```

ソースコード 15 1_5.py

```

1 #2019/04/08
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4
5 import sys
6
7 def decode(text, offset):
8     length = len(text)
9     #if too short
10     if length < offset + 3:
11         return ''
12     peptide = ''
13     CODON = {'TTT':'Phe', 'TTC':'Phe', 'TTA':'Leu', 'TTG':'Leu',\
14             'TCT':'Ser', 'TCC':'Ser', 'TCA':'Ser', 'TCG':'Ser',\
15             'TAT':'Tyr', 'TAC':'Tyr', 'TAA':'_', 'TAG':'_',\
16             'TGT':'Cys', 'TGC':'Cys', 'TGA':'_', 'TGG':'Trp',\
17             'CTT':'Leu', 'CTC':'Leu', 'CTA':'Leu', 'CTG':'Leu',\
18             'CCT':'Pro', 'CCC':'Pro', 'CCA':'Pro', 'CCG':'Pro',\
19             'CAT':'His', 'CAC':'His', 'CAA':'Gln', 'CAG':'Gln',\
20             'CGT':'Arg', 'CGC':'Arg', 'CGA':'Arg', 'CGG':'Arg',\
21             'ATT':'Ile', 'ATC':'Ile', 'ATA':'Ile', 'ATG':'Met',\
22             'ACT':'Thr', 'ACC':'Thr', 'ACA':'Thr', 'ACG':'Thr',\
23             'AAT':'Asn', 'AAC':'Asn', 'AAA':'Lys', 'AAG':'Lys',\
24             'AGT':'Ser', 'AGC':'Ser', 'AGA':'Arg', 'AGG':'Arg',\
25             'GTT':'Val', 'GTC':'Val', 'GTA':'Val', 'GTG':'Val',\
26             'GCT':'Ala', 'GCC':'Ala', 'GCA':'Ala', 'GCG':'Ala',\
27             'GAT':'Asp', 'GAC':'Asp', 'GAA':'Glu', 'GAG':'Glu',\
28             'GGT':'Gly', 'GGC':'Gly', 'GGA':'Gly', 'GGG':'Gly'}
29     start = offset
30     end = start + 3
31     try:
32         #Skip to Met
33         while end <= length:
34             if CODON[text[start:end]] == 'Met':
35                 break
36             start = end
37             end = start + 3
38         #Decode
39         while end <= length:
40             peptide = peptide + ' ' + CODON[text[start:end]]
41             if peptide[len(peptide)-1] == '_':
42                 break
43             start = end
44             end = start + 3
45         return peptide.strip()
46     except (KeyError):
47         sys.stderr.write('Unknown CODON Error\n')

```

```
48 |         return ''
```

ソースコード 16 decode.py

実行した結果が以下の通りである.

```
1 PEPTIDE 1
2 Met Ile Val Met Asn Ser Asn Cys Cys Leu Cys Arg Pro Thr Arg Phe Leu Thr Ser Leu Ser Tyr His Phe
3   Leu Leu Ser Tyr Leu Leu Ser Lys Cys Ile Gln Met Lys Gly Cys Gly Glu Cys _
4 PEPTIDE 2
5 Met Pro Arg Glu Ile Ser Arg Ser Ser Val Pro Cys _
6 PEPTIDE 3
7 Met Pro _
8 PEPTIDE 4
9 Met Leu Arg Thr Leu Leu Leu Ile Val _
10 PEPTIDE 5
11 Met Gln Asn Lys Phe Ile Arg His _
12 PEPTIDE 6
13 Met Gly Arg Lys Asp Lys Ala Ala Ile _
```

ソースコード 17 実行結果

2 タンパク質構造情報処理

2.1 PyMOL ソフトウェアを利用し、ヒトのヘモグロビンの構造をチェーンごとに異なる色で表示し、4 量体であることを確認せよ。また、A チェインだけを表示し、タンパク質鎖を cartoon 表示して 2 次構造に従って色付けし、結合するヘムを stick 表示して原子ごとに色分けせよ。

チェーンごとに色分けしたヘモグロビンが画像 2 である。

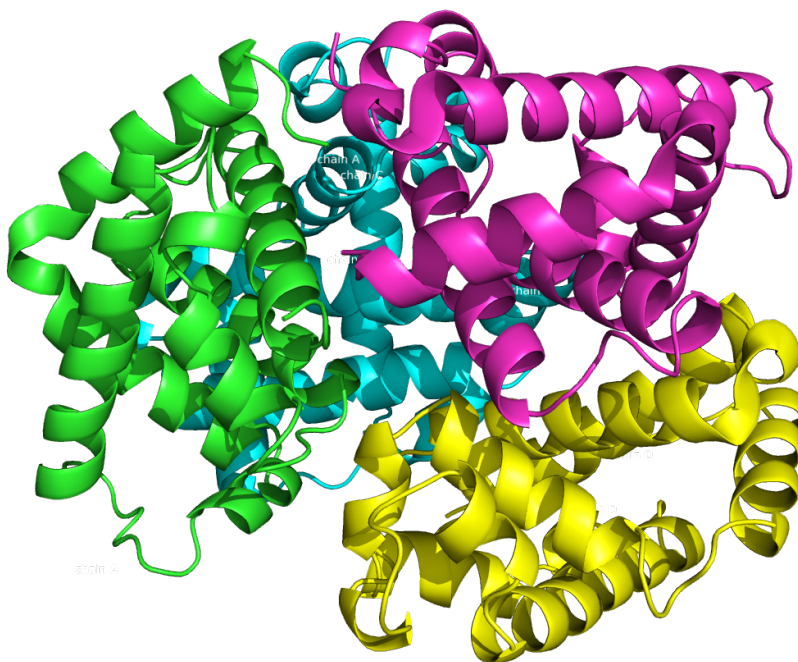


図 2 ヘモグロビン

また、A チェインとヘムを表示したものが画像 3 である。

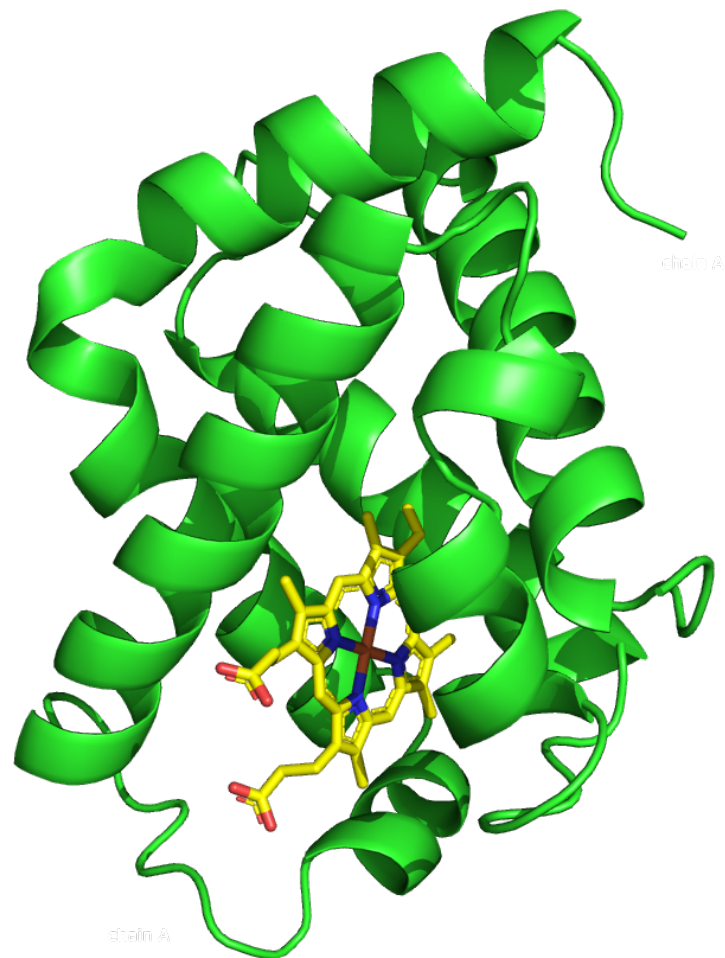


図 3 A チェインとヘム

2.2 PDB ファイル名とチェーン名を引数にして、その回転半径を計算するプログラムを作成せよ。

2.3 上記のプログラムの結果を利用し、PyMOL で重心から回転半径の範囲内にある原子を赤で、範囲外にある原子を青で色付けせよ。

この 2 つは一緒に行い、以下のプログラムを作成した。実行環境は Anaconda3 の Python 2.7 環境である。

```
1 #2019/04/17
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import math
7 import sys
8 import pymol
```



```

9
10 #file open
11 args = sys.argv
12 if len(args) != 3:
13     sys.stderr.write('Usage: python 2_7.py ./1BUW.pdb A\n')
14     sys.exit()
15 path = args[1]
16 try:
17     file = open(path)
18 except:
19     sys.stderr.write('Error in opening file.\n')
20     sys.exit()
21 lines = file.readlines()
22
23 WEIGHT = \
24 { 'H':1.008,\
25   'LI':6.941, 'BE':9.012, 'B':10.81, 'C':12.01, 'N':14.01, 'O':16.00, 'F':19.00,\
26   'NA':22.99, 'MG':24.31, 'AL':26.98, 'SI':28.09, 'P':30.97, 'S':32.07, 'CL':35.45,}
27 total_weight = 0
28 x_weight = 0
29 y_weight = 0
30 z_weight = 0
31 atom_data = []
32 for line in lines:
33     if line[0:6] == 'ATOM' and line[21] == args[2]:
34         #Parse Data
35         atom_number = int(line[6:11])
36         atom_name = line[12:16].strip()
37         #identifier = line[16].strip()
38         #res_name = line[17:20].strip()
39         chain = line[21].strip()
40         res_number = int(line[22:26])
41         #res_code = line[26].strip()
42         x = float(line[30:38])
43         y = float(line[38:46])
44         z = float(line[46:54])
45         #occupy = float(line[54:60])
46         #tmp = line[60:66].strip()
47         elem_symbol = line[76:78].strip()
48         #charge = line[78:80].strip()
49
50         #Calculate Center of Gravity
51         total_weight = total_weight + WEIGHT[elem_symbol]
52         x_weight = x_weight + x * WEIGHT[elem_symbol]
53         y_weight = y_weight + y * WEIGHT[elem_symbol]
54         z_weight = z_weight + z * WEIGHT[elem_symbol]
55
56         #Register Atom Data
57         atom_data.append([x, y, z, atom_number, atom_name, elem_symbol, chain, res_number])
58 #IF no Data
59 if total_weight == 0:
60     sys.stderr.write('Chain '+args[2]+' does not exist.\n')
61     sys.exit()
62
63 #Calculate Center of Gravity
64 x_balance = x_weight / total_weight
65 y_balance = y_weight / total_weight
66 z_balance = z_weight / total_weight
67
68 #Calculate circle r
69 r = 0
70 for data in atom_data:
71     r = r + math.sqrt((data[0]-x_balance)**2+(data[1]-y_balance)**2+(data[2]-z_balance)**2)
72 r = r / len(atom_data)
73 print('Center: ('+str(x_balance)+' , '+str(y_balance)+' , '+str(z_balance)+')')
74 print('Radius: '+str(r))
75
76 #identify color
77 red = []
78 blue = []
79 for data in atom_data:
80     if math.sqrt((data[0]-x_balance)**2+(data[1]-y_balance)**2+(data[2]-z_balance)**2) > r:
81         blue.append(data)
82     else:
83         red.append(data)
84
85 #Launch PyMOL
86 pymol.finish_launching()
87 pymol.cmd.load(path)
88 pymol.cmd.do("hide everything")
89 pymol.cmd.do("show stick, chain "+args[2])

```

```

90 #set color
91 for data in red:
92     pymol.cmd.do("color red, (name "+data[4]+" and resi "+str(data[7])+
93     " and chain "+data[6]+")")
94 for data in blue:
95     pymol.cmd.do("color blue, (name "+data[4]+" and resi "+str(data[7])+
96     " and chain "+data[6]+")")
97 file.close()

```

ソースコード 18 2_7.py

このプログラムを次のように実行する。

```

1 python 2_7.py ./1BUW.pdb A

```

ソースコード 19 実行方法

結果は次の通りである。

```

1 Center: (47.0059686237, 33.4827138428, 35.7021762506)
2 Radius: 14.0781227416

```

ソースコード 20 実行結果

また、描画結果として以下の画像を得た。

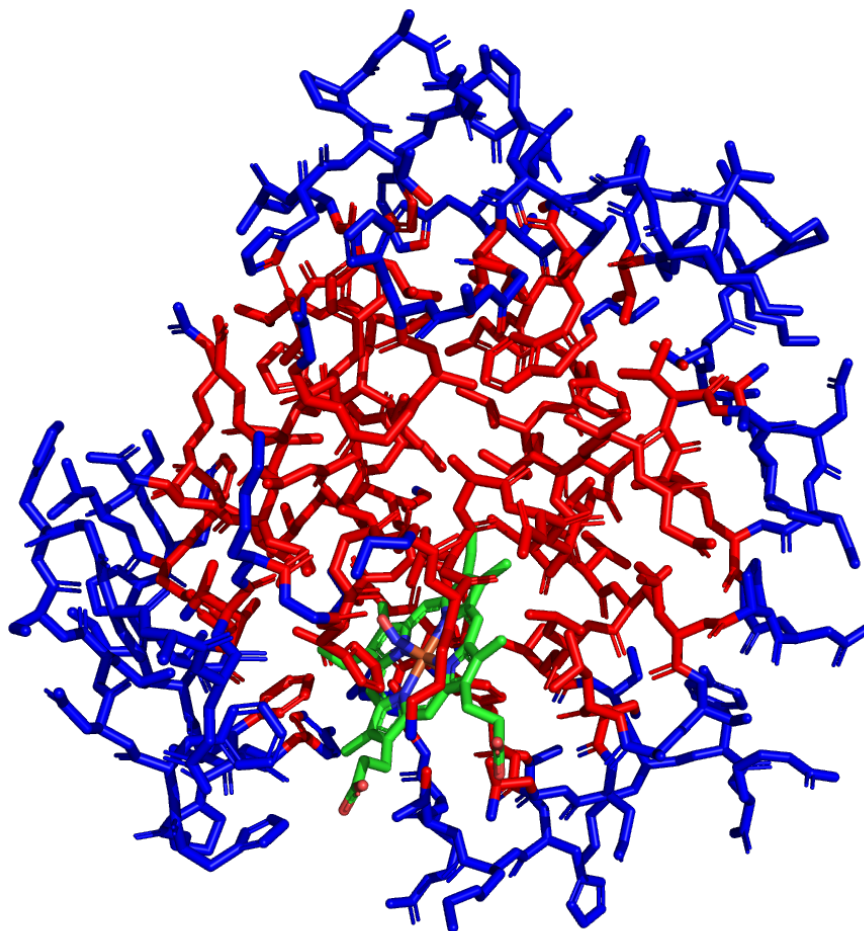


図 4 回転半径

3 機械学習

3.1 Python の scikit-learn の SVM で ionosphere のデータの 10-fold Cross Validation を実施せよ。予測性能として precision, recall, MCC, ROC 曲線の AUC 値 (AUROC), F-score を求めよ。カーネルは RBF カーネルとし、パラメータは適宜定めよ。

作成したプログラムを 21 に示す。Anaconda3 の Python2.7 環境で動作する。

```
1 #2019/04/16
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7 from sklearn.svm import SVC
8 from sklearn.metrics import precision_score, recall_score, f1_score, matthews_corrcoef,
9   roc_auc_score
10 from load_ionosphere import load_ionosphere
11 from sklearn.model_selection import *
12
13 #TP, TN : True
14 #FN : actually positive, but expected negative
15 #FP : actually negative, but expected positive
16 #Precision : TP / TP + FP
17   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.precision_score.html#
18   sklearn.metrics.precision_score
19 #Recall : TP / TP + FN
20   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.recall_score.html#
21   sklearn.metrics.recall_score
22 #F-measure : 2 * Recall * Precision / Recall + Precision
23   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html#sklearn.
24   metrics.f1_score
25 #MCC : Matthews Correlation Coefficient
26   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.matthews_corrcoef.html#
27   sklearn.metrics.matthews_corrcoef
28 #AUROC :
29   # https://scikit-learn.org/stable/modules/generated/sklearn.metrics.roc_auc_score.html#
30   sklearn.metrics.roc_auc_score
31
32 svc = SVC(C=1.0, kernel='rbf', gamma=0.2)
33 kf = KFold(n_splits=10, shuffle=True)
34 [target, parameters] = load_ionosphere()
35 counter = 0
36 print('Set\tPreci.\tRecall\tMCC\tAUROC\tF-score')
37 for train_index, test_index in kf.split(parameters, target):
38     counter = counter + 1
39     train_param = []
40     train_target = []
41     test_param = []
42     test_target = []
43     for i in range(0, len(train_index)):
44         train_param.append(parameters[train_index[i]])
45         train_target.append(target[train_index[i]])
46     for i in range(0, len(test_index)):
47         test_param.append(parameters[test_index[i]])
48         test_target.append(target[test_index[i]])
49     sys.stdout.write(str(counter)+'\t')
50     learn = svc.fit(train_param, train_target)
51     print(str(precision_score(test_target, learn.predict(test_param)))[0:7]+' \t' \
52           +str(recall_score(test_target, learn.predict(test_param)))[0:7]+' \t' \
53           +str(matthews_corrcoef(test_target, learn.predict(test_param)))[0:7]+' \t' \
54           +str(roc_auc_score(test_target, learn.predict(test_param)))[0:7]+' \t' \
55           +str(f1_score(test_target, learn.predict(test_param)))[0:7])
```

ソースコード 21 3.9.py

また、ionosphere.scale の読み込みにあたって、欠番のパラメータを 0 で補完するなどの処理を行うため、読み込み部分は別のプログラム 22 として分離した。

```
1 #2019/04/16
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
```

```

4 #For Python 2.7
5
6 import sys
7
8 def load_ionosphere():
9     #file open
10    try:
11        file = open('./ionosphere.scale')
12    except:
13        sys.stderr.write('Error in opening file.\n')
14        sys.exit()
15
16    lines = file.readlines()
17    target = []
18    parameters = []
19    for line in lines:
20        set = line.split()
21        #target is 1 or 0
22        target.append(int(set.pop(0))+1//2)
23        datas = {}
24        params = []
25        for data in set:
26            datas[int(data.split(':')[0])] = float(data.split(':')[1])
27        for i in range(0,34):
28            try:
29                params.append(datas[i+1])
30            except KeyError:
31                params.append(0)
32            parameters.append(params)
33    file.close()
34    return [target, parameters]

```

ソースコード 22 load_ionosphere.py

実行した結果を 23 に示す.

| Set | Preci. | Recall | MCC | AUROC | F-score |
|-----|---------|---------|---------|---------|---------|
| 1 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| 2 | 1.0 | 0.95454 | 0.94146 | 0.97727 | 0.97674 |
| 3 | 0.9 | 1.0 | 0.89113 | 0.94117 | 0.94736 |
| 4 | 0.84 | 0.95454 | 0.69186 | 0.82342 | 0.89361 |
| 5 | 0.92 | 1.0 | 0.87559 | 0.91666 | 0.95833 |
| 6 | 1.0 | 1.0 | 1.0 | 1.0 | 1.0 |
| 7 | 0.91666 | 0.95652 | 0.80760 | 0.89492 | 0.93617 |
| 8 | 1.0 | 0.9375 | 0.75 | 0.96875 | 0.96774 |
| 9 | 1.0 | 0.95238 | 0.94280 | 0.97619 | 0.97560 |
| 10 | 0.90476 | 0.95 | 0.82495 | 0.90833 | 0.92682 |

ソースコード 23 実行結果

3.2 ionosphere データにおいて、より良い RBF カーネルパラメータ γ とコストパラメータ C の値を探索せよ. 評価方法は 10-fold Cross Validation とし、評価基準は AUROC と F-score の 2 通りを試すこと.

作成したプログラムを 24 に示す. Anaconda3 の Python2.7 環境で動作する.

```

1 #2019/04/16
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 2.7
5
6 import sys
7 from sklearn.svm import SVC
8 from sklearn.metrics import precision_score, recall_score, f1_score, matthews_corrcoef,
9   roc_auc_score
10 from load_ionosphere import load_ionosphere
11 from sklearn.model_selection import *
12
13 #Help https://scikit-learn.org/stable/modules/generated/sklearn.model_selection.GridSearchCV.html
14 #Help https://scikit-learn.org/stable/modules/model_evaluation.html#scoring-parameter
15
16 [target, parameters] = load_ionosphere()

```

```

16 param_grid = {'C':[0.001, 0.01, 0.1, 1, 10, 100, 1000], 'gamma':[0.001, 0.01, 0.1, 1, 10, 100,
17                  1000]}
18 print('Evaluate with AUROC')
19 search_auroc = GridSearchCV(estimator=SVC(), param_grid=param_grid, scoring='roc_auc', cv=10,
20                             iid=False).fit(parameters, target)
21 print(search_auroc.best_params_)
22 print(search_auroc.best_score_)
23 print('Evaluate with F-score')
24 search_f = GridSearchCV(estimator=SVC(), param_grid=param_grid, scoring='f1', cv=10, iid=False).
25 fit(parameters, target)
26 print(search_f.best_params_)
27 print(search_f.best_score_)

```

ソースコード 24 3_10.py

実行した結果を 25 に示す.

```

1 Evaluate with AUROC
2 {'C': 10, 'gamma': 1}
3 0.9860481909394953
4 Evaluate with F-score
5 {'C': 10, 'gamma': 0.1}
6 0.9616653503831382

```

ソースコード 25 実行結果

4 創薬情報処理・機械学習

4.1 (準備 1) Python に RDKit をインストールし, 化合物ファイル (SDF ファイル) を読み込んで構造式が出力できることを確認せよ.

省略する.

4.2 (準備 2) 論文 Leung SSF, et al. J Chem Inf Model 56: 924-020, 2016. の Supporting Information より, 3D SDF ファイル (TXT) と PDF ファイルをダウンロードせよ. Table S1 の実験値「RRCK Log P_{app} 」の値をパース (転記) し, CSV ファイル等で準備せよ.

省略する.

4.3 (準備 3) 所望の化合物に対し, RDKit の ECFP4 fingerprint を計算できるようにせよ.

calc_fgprint 関数を作成した。コードは 26 の通りである。

```

1 #2019/04/15
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 3.7
5
6 #User rdkit-env
7
8 from rdkit import Chem
9 from rdkit.Chem import AllChem
10
11 def calc_fgprint(mol):
12     data = AllChem.GetMorganFingerprintAsBitVect(mol, 2, nBits = 1024)
13     set = []
14     for i in range(0,1024):
15         if data.GetBit(i):
16             set.append(1)
17         else:

```

```

18         set.append(0)
19     return set

```

ソースコード 26 calc_fgprint.py

4.4 RRCK Log P_{app} を目的関数, ECFP4 fingerprint を特徴ベクトルとして, Data Set 3(医薬品 104 化合物) に関して回帰 (10-fold Cross Validation) を行う機械学習プログラムを作成せよ. 学習器は Support Vector Regression とし, カーネルは RBF カーネルとすること.

スクリプト 27 を作成した。

```

1  #2019/04/15
2  #Yuki Tsushima
3  #tsushima@bi.c.titech.ac.jp
4  #For Python 3.7
5
6  #User rdkit-env
7
8  import csv
9  import sys
10 from calc_fgprint import calc_fgprint
11 from rdkit import rdBase, Chem
12 from rdkit.Chem import AllChem, Draw
13 from rdkit.Chem.Draw import rdMolDraw2D
14
15 from sklearn import *
16 from sklearn.svm import SVR
17 from sklearn.model_selection import *
18
19 #Load Mols
20 suppl = Chem.SDMolSupplier('./ci6b00005_si_002.txt', removeHs=False) #Supplier
21 mols = [x for x in suppl if x is not None]
22 mols_set3 = mols[106:210]
23 fingerprint = []
24 for data in mols_set3:
25     fingerprint.append(calc_fgprint(data))
26
27 #Load RRCKs
28 csv_file = open('./dataset3.csv', 'r')
29 f = csv.reader(csv_file)
30 RRCK_set3 = []
31 for row in f:
32     RRCK_set3.append(float(row[0]))
33
34 svr = SVR(kernel='rbf', C=100, gamma=0.01)
35
36 score = cross_val_score(estimator=svr, X=fingerprint, y=RRCK_set3, cv=10)
37 print(sum(score)/len(score))

```

ソースコード 27 4_14.py

実行結果は 28 の通りである。

```

1  0.2630136419201868

```

ソースコード 28 実行結果

4.5 10-fold Cross Validation によって, RRCK Log P_{app} の予測値との平均 2 乗誤差 (RMSE) が最も小さくなるパラメータを探索せよ. またそのときの RMSE と R^2 値を求めよ.

スクリプト 29 を作成した。

```

1  #2019/04/18

```

```

2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 3.7
5
6 #User rdkit-env
7
8 import csv
9 import sys
10 from calc_fgprint import calc_fgprint
11 from rdkit import rdBase, Chem
12
13 from sklearn import svm, model_selection, metrics
14 from sklearn.svm import SVR
15 from sklearn.model_selection import GridSearchCV
16 from sklearn.metrics import mean_squared_error, r2_score
17 import numpy as np
18
19 #Load Mols
20 suppl = Chem.SDMolSupplier('./ci6b00005_si_002.txt', removeHs=False) #Supplier
21 mols = [x for x in suppl if x is not None]
22 mols_set3 = mols[106:210]
23 fingerprint = []
24 for data in mols_set3:
25     fingerprint.append(calc_fgprint(data))
26
27 #Load RRCKs
28 csv_file = open('./dataset3.csv', 'r')
29 f = csv.reader(csv_file)
30 RRCK_set3 = []
31 for row in f:
32     RRCK_set3.append(float(row[0]))
33
34 param_grid = {'C':[0.001, 0.01, 0.1, 1, 10, 100, 1000], 'gamma':[0.001, 0.01, 0.1, 1, 10, 100, 1000]}
35 def rmse(y, y_pred):
36     return np.sqrt(mean_squared_error(y, y_pred))
37 search = GridSearchCV(estimator=SVR(), param_grid=param_grid, scoring='neg_mean_squared_error',
38 cv=10, iid=False).fit(fingerprint, RRCK_set3)
39 print('R2_score : {}'.format(r2_score(RRCK_set3, search.predict(fingerprint))))
40 print('RMSE : {}'.format(rmse(RRCK_set3, search.predict(fingerprint))))
41 print('best_params : {}'.format(search.best_params_))

```

ソースコード 29 4_15.py

実行結果は 30 の通りである。

```

1 R2_score : 0.9747505807591622
2 RMSE : 0.10158867737298335
3 best_params : {'C': 100, 'gamma': 0.01}

```

ソースコード 30 実行結果

4.6 (15) で決めたパラメータの学習器で、Data Set 1(環状ペプチド 7 化合物)、Data Set 4(環状ペプチド 16 化合物)、Data Set 8(環状ペプチド 22 化合物) の RRCK Log P_{app} の予測値を求めよ。それぞれの Data Set ごとに RMSE と R^2 値を求め、予測値と実験値の散布図を描け。

スクリプト 31 を作成した。

```

1 #2019/04/18
2 #Yuki Tsushima
3 #tsushima@bi.c.titech.ac.jp
4 #For Python 3.7
5
6 #User rdkit-env
7
8 import codecs
9 import csv
10 import sys
11 from calc_fgprint import calc_fgprint
12 from rdkit import rdBase, Chem
13

```

```

14 from sklearn import svm, model_selection, metrics
15 from sklearn.svm import SVR
16 from sklearn.model_selection import GridSearchCV
17 from sklearn.metrics import mean_squared_error, r2_score
18 import numpy as np
19
20 #Load Mols
21 suppl1 = Chem.SDMolSupplier('./dataset1.sdf', removeHs=False) #Supplier
22 suppl3 = Chem.SDMolSupplier('./ci6b00005_si_002.txt', removeHs=False) #Supplier
23 suppl4 = Chem.SDMolSupplier('./dataset4.sdf', removeHs=False) #Supplier
24 suppl8 = Chem.SDMolSupplier('./dataset8.sdf', removeHs=False) #Supplier
25 mols = [x for x in suppl3 if x is not None]
26 mols_set3 = mols[106:210]
27 mols_set1 = [x for x in suppl1 if x is not None]
28 mols_set4 = [x for x in suppl4 if x is not None]
29 mols_set8 = [x for x in suppl8 if x is not None]
30 finger1 = []
31 finger3 = []
32 finger4 = []
33 finger8 = []
34 for data in mols_set1:
35     finger1.append(calc_fgprint(data))
36 for data in mols_set3:
37     finger3.append(calc_fgprint(data))
38 for data in mols_set4:
39     finger4.append(calc_fgprint(data))
40 for data in mols_set8:
41     finger8.append(calc_fgprint(data))
42
43 #Load RRCKs
44 csv1 = open('./dataset1.csv', 'r')
45 f1 = csv.reader(csv1)
46 RRCK_set1 = []
47 for row in f1:
48     RRCK_set1.append(float(row[0]))
49 csv3 = open('./dataset3.csv', 'r')
50 f3 = csv.reader(csv3)
51 RRCK_set3 = []
52 for row in f3:
53     RRCK_set3.append(float(row[0]))
54 csv4 = open('./dataset4.csv', 'r')
55 f4 = csv.reader(csv4)
56 RRCK_set4 = []
57 for row in f4:
58     RRCK_set4.append(float(row[0]))
59 csv8 = open('./dataset8.csv', 'r')
60 f8 = csv.reader(csv8)
61 RRCK_set8 = []
62 for row in f8:
63     RRCK_set8.append(float(row[0]))
64
65 #Learn
66 svr = SVR(kernel='rbf', C=100, gamma=0.01).fit(finger3, RRCK_set3)
67 def rmse(y, y_pred):
68     return np.sqrt(mean_squared_error(y, y_pred))
69
70 #Predict and output
71 pre1 = svr.predict(finger1)
72 with codecs.open('predict1.csv', 'w', 'utf-8') as f:
73     for i in range(0, len(pre1)):
74         f.write(str(RRCK_set1[i])+' '+str(pre1[i])+'\n')
75 print("Dataset 1")
76 print('R2_score : {}'.format(r2_score(RRCK_set1, pre1)))
77 print('RMSE : {}'.format(rmse(RRCK_set1, pre1)))
78 pre3 = svr.predict(finger3)
79 with codecs.open('predict3.csv', 'w', 'utf-8') as f:
80     for i in range(0, len(pre3)):
81         f.write(str(RRCK_set3[i])+' '+str(pre3[i])+'\n')
82 print("Dataset 3")
83 print('R2_score : {}'.format(r2_score(RRCK_set3, pre3)))
84 print('RMSE : {}'.format(rmse(RRCK_set3, pre3)))
85 pre4 = svr.predict(finger4)
86 with codecs.open('predict4.csv', 'w', 'utf-8') as f:
87     for i in range(0, len(pre4)):
88         f.write(str(RRCK_set4[i])+' '+str(pre4[i])+'\n')
89 print("Dataset 4")
90 print('R2_score : {}'.format(r2_score(RRCK_set4, pre4)))
91 print('RMSE : {}'.format(rmse(RRCK_set4, pre4)))
92 pre8 = svr.predict(finger8)
93 with codecs.open('predict8.csv', 'w', 'utf-8') as f:
94     for i in range(0, len(pre8)):

```



```

95     f.write(str(RRCK_set8[i])+','+str(pre8[i]))+'\n')
96 print("Dataset 8")
97 print('R2_score : {}'.format(r2_score(RRCK_set8, pre8)))
98 print('RMSE : {}'.format(rmse(RRCK_set8, pre8)))

```

ソースコード 31 4_16.py

実行結果は 32 の通りである。

```

1 Dataset 1
2 R2_score : -2.314801681170682
3 RMSE : 0.6643257009105027
4 Dataset 3
5 R2_score : 0.9747505807591622
6 RMSE : 0.10158867737298335
7 Dataset 4
8 R2_score : -2.1150984682038985
9 RMSE : 0.9198261864774264
10 Dataset 8
11 R2_score : 0.12606435549538175
12 RMSE : 0.4296164950959734

```

ソースコード 32 実行結果

また、実験値と予測値の散布図は出力したファイルから Gnuplot で描画した。図 5、6、7、8 を得た。

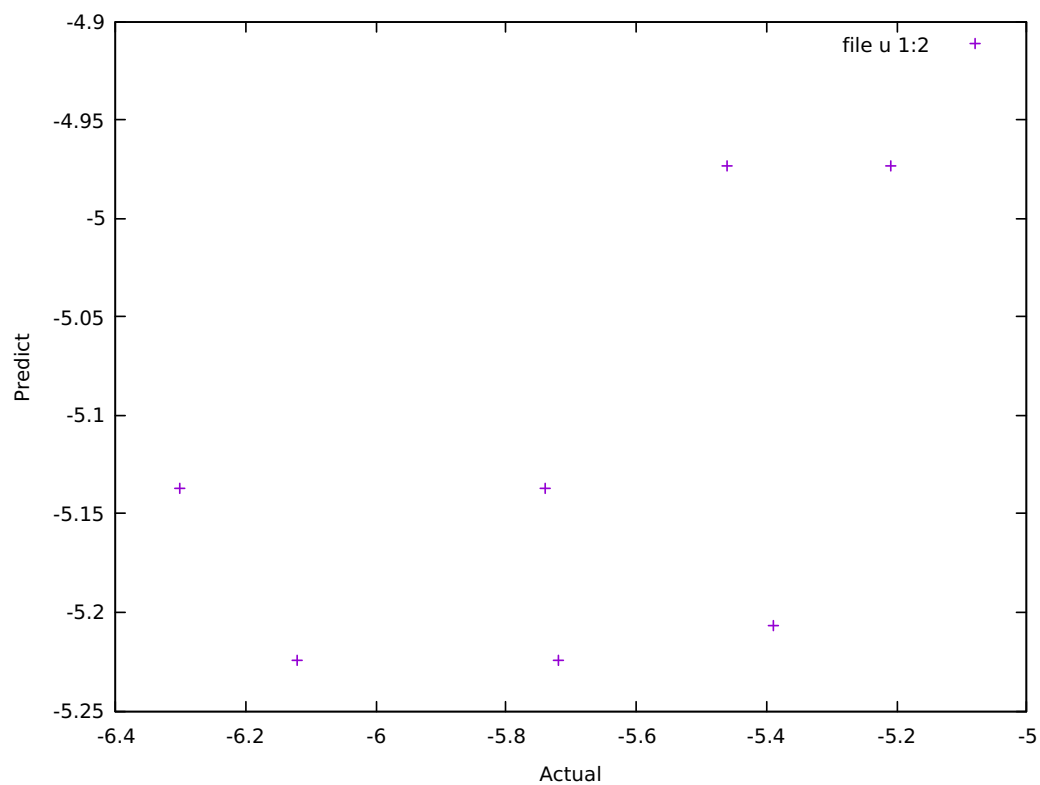


図 5 Data Set 1

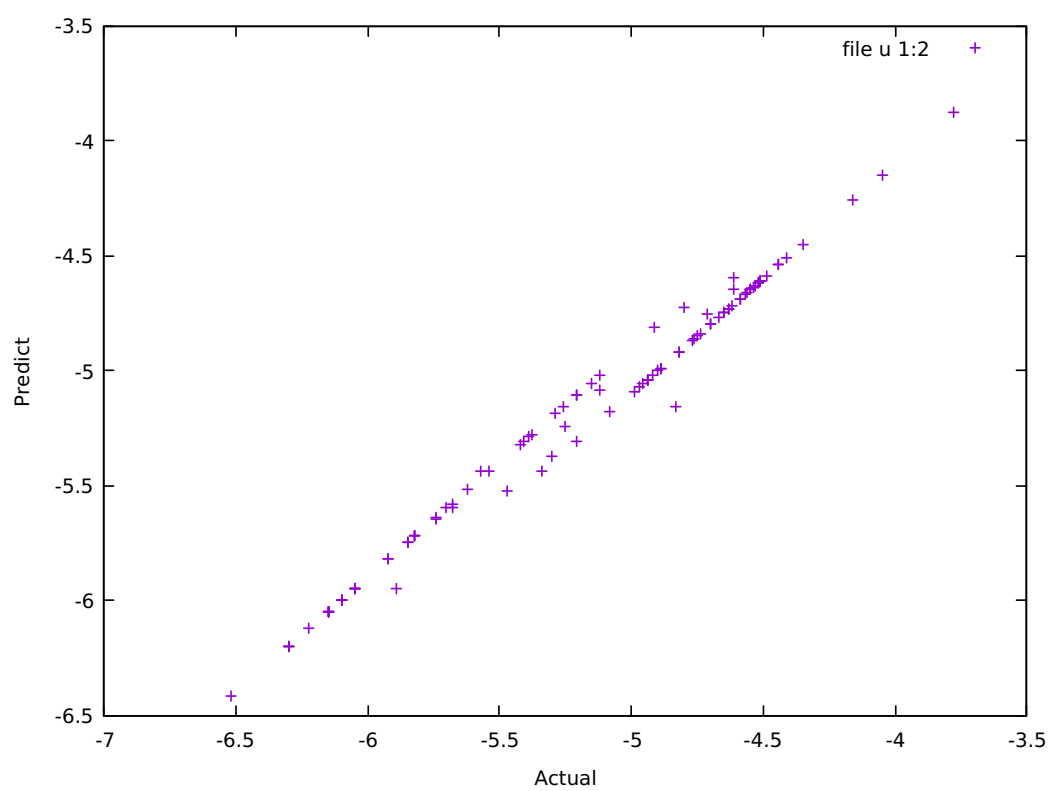


图 6 Data Set 3

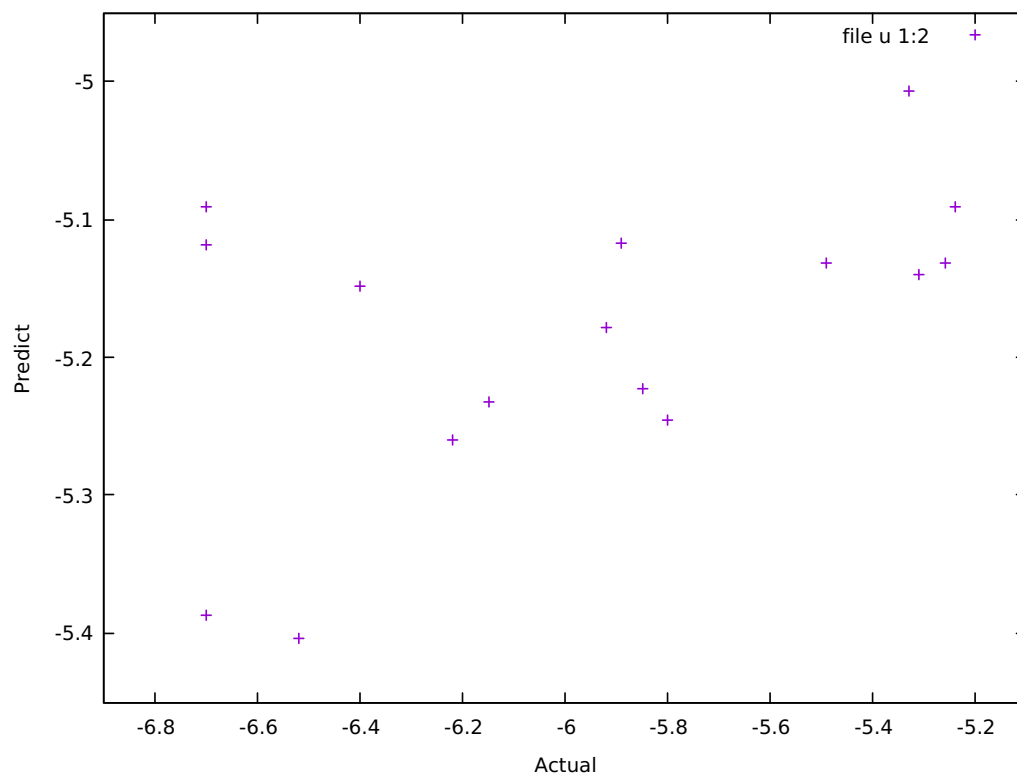


图 7 Data Set 4

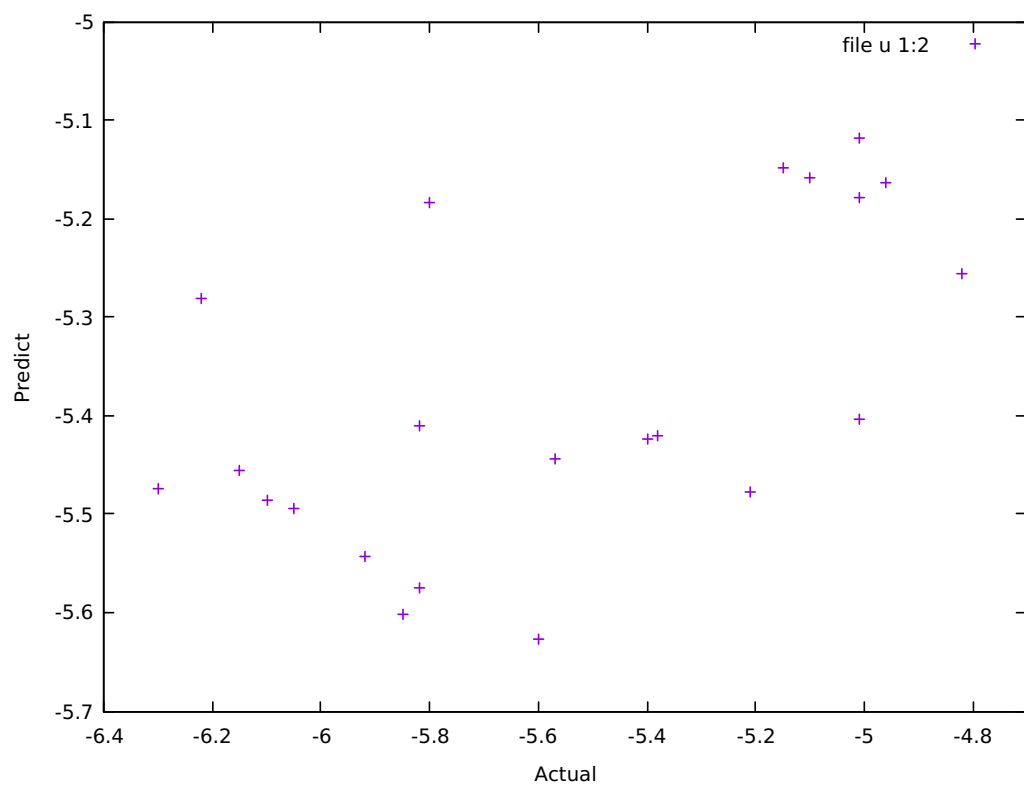


图 8 Data Set 8