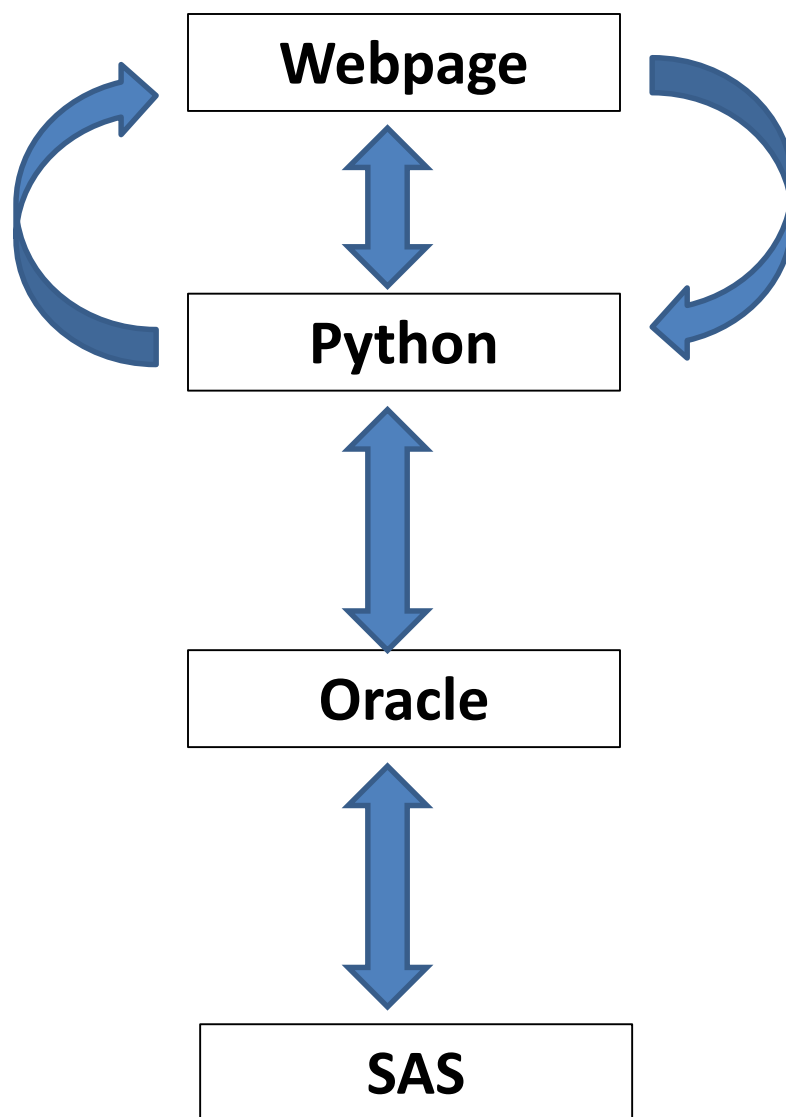# STSCI 4060

# Lecture File 11

**Xiaolong Yang**
**(xy44@cornell.edu)**

# **Python, Database-Driven Web Pages, and SAS**

# Python, Oracle, Webpage and SAS Integration

# Prepare/Gather Data of Interest

Where to get data?

- Data that people already collected and managed in a format ready for use

- Data that have been collected but need treatment

- Data that are somewhat all over different places, e.g., internet

- Data to be collected by you or someone through various experiments

# **Where and How to Store your Data?**

- Your personal media, e.g., disks, flash drives, etc.

- Your personal computer (flat files, MS Excel files, etc.)

- Databases of different levels

  - Local databases

  - Global databases connected via Internet

# Our Scenario

- Collect data of interest from internet

- Study the data collected

- Treat/transform the data with Python programs

- Store the transformed data in an Oracle database through a network using Python

- Access the database through a network using Python, which is an example of database-driven, Python-powered dynamic webpage

- Analyze the data stored in the database with SAS, etc.

# Gather Our Data of Interest

- Honey bee proteins that have been discovered and studied.
- Available from the Internet, especially from NCBI (http://www.ncbi.nlm.nih.gov/) , for example

```
>gi|163311124|pdb|3BJH|A Chain A, Soft-Sad Crystal Structure Of A Pheromone Bi
Protein From The Honeybee Apis Mellifera L.
APDWVPPEVFDLVAEDKARCMSEHGTTQAQIDDVDKGNLVNEPSITCYMYCLLEAFSLVDDEANVDEDIM
LGLLPDQLQERAQSVMGKCLPTSGSDNCNKIYNLAKCVQESAPDVWFVI
```

☐ vitellogenin precursor [Apis mellifera]
3. NCBI Reference Sequence: NP_001011578.1
GenPept    Graphics    Related Sequences    Identical Proteins

Previous record

```
>gi|58585104|ref|NP_001011578.1| vitellogenin precursor [Apis mellifera]
MLLLLTLLLFAGTVAADFQHNWQVGNEYTYLVRSRTLTSLGDLSDVHTGILIKALLTVQAKDSNVLAAKV
WNGQYARVQQSMPDGWETEISDQMLELRDLPISGKPFQIRMKHGLIRDLIVDRDVPTWEVNILKSIVGQL
QVDTQGENAVKVNSVQVPTDDEPYASFKAMEDSVGGKCEVLYDIAPLSDFVIHRSPELVPMPTLKGDGRH
MEVIKIKNFDNCDQRINYHFGMTDNSRLEPGTNKNGKFFSRSSTSRIVISESLKHFTIQSSVTTSKMMVS
PRLYDRQNGLVLSRMNLTLAKMEKTSKPLPMVDNPESTGNLVYIYNNPFSDVEERRVSKTAMNSNQIVSD
NSLSSSEEKLKQDILNLRTDISSSSSSISSSEENDFWQPKPTLEDAPQNSLLPNFVGYKGKHIGKSGKVD
VINAAKELIFQIANELEDASNIPVHATLEKFMILCNLMRTMNRKQISELESNMQISPNELKPNDKSQVIK
QNTWTVFRDAITQTGTGPAFLTIKEWIERGTTKSMEAANIMSKLPKTVRTPTDSYIRSFFELLQNPKVSN
EQFLNTAATLSFCEMIHNAQVNKRSIHNNYPVHTFGRLTSKHDNSLYDEYIPFLERELRKAHQEKDSPRI
QTYIMALGMIGEPKILSVFEPYLEGKQQMTVFQRTLMVGSLGKLTETNPKLARSVLYKIYLNTMESHEVR
CTAVFLLMKTNPPLSMLQRMAEFTKLDTNRQVNSAVKSTIQSLMKLKSPEWKDLAKKARSVNHLLTHHEY
DYELSRGYIDEKILENQNIITHMILNYVGSEDSVIPRILYLTWYSSNGDIKVPSTKVLAMISSVKSFMEL
```

# The Honey Bee Protein Data Set

# Inspect The Honey Bee Protein Data Set

- An entry (protein) always starts with "**>gi|**" (in our case)
- It starts with a short description about the protein, but the length varies
- There is a substring "**[*Apis mellifera*]** " right before the actual protein sequence
- The sequence has a **carriage return** at the end of each line of the sequence
- Our interest:
  - \* Only interested in one species of honey bees, *Apis mellifera*
  - \* The gi number (which can be used as a protein ID)
  - \* The continuous protein sequence
  - \* The relative frequency of the amino acids in each protein

# Our Plan to Treat the Data Set

- Use Python as the programming tool

- Submit this data file to the database (Oracle) through a specific web page

- Read in the content as a string

- Change content into a continuous string

- Split the string into a list of bee proteins that only belong to *Apis mellifera*

- Extract the gi number and the sequence of each protein

- Calculate the relative frequencies of the 20 amino acids

- Store the data/information in the Oracle database through a network

# Design a Web Page to Upload the Data

**Browser View →**



**KompoZer Design View →**

# Python Coding: the General Algorithm

- Code two CGI (common gateway interface) files, one for reading and processing data and one for displaying database contents.
- Use Python functions to achieve the goal: main(), processInput(), …
- The main() function receives the file name from the web page and pass it to the processInput() function and finally display the results on a web page.
- The processInput() function creates a new database table in Oracle, reads and processes the data as planned earlier. More details on database handling:
    Use special features (e.g., ">gi|", "[Apis mellifera]") of the
    dataset to extract the gi numbers and protein sequences:
        *   Use bind variables to populate the database table.
        *   Write other functions  when necessary, e.g., fileToStr() and makePage().
- Design two webpage templates, one to confirm uploading data to the database and one to display some contents retrieved from the database.

# The 1st CGI File

```
1   '''
2   This program processes the data file uploaded by the
3   user through a network and saves the processed data to an
4   Oracle database.
5
6   ******** Coded by Xiaolong Yang @ Cornell University ********
7   '''
8
9   import cgi
10  import cx_Oracle
11
12  def main(): # NEW
13      form=cgi.FieldStorage() #cgi script line
14      theStr=form.getfirst('theList','')
15      contents = processInput(theStr)
16      print contents
17
18  def processInput(theFile):
19      '''
20      This function reads the bee protein data from a raw data file and
21      extract the gi number values and the protein sequences. Then, the
22      relative frequencies of the amino acids are calculated for each protain.
23      A new Oracle database table is created using Python and is populated with
24      the data created.
25      '''
26
27      con = cx_Oracle.connect('python/welcome')
28      cur=con.cursor()
```

# The 1ˢᵗ CGI File (cont'd)

query_arraysize.sql ☒ | bind_insert.sql ☒ | Data_output_from_DB_dynamic.py ☒ | Data_Input_dynamic.cgi ☒ | notepatplus.html ☒

```
28        cur=con.cursor()
29        cur.execute('drop table beeProteins')
30        cur.execute('''create table beeProteins (
31                        gi varchar2(10),
32                        sequence clob,
33                        freq_A number,
34                        freq_R number,
35                        freq_N number,
36                        freq_D number,
37                        freq_C number,
38                        freq_Q number,
39                        freq_E number,
40                        freq_G number,
41                        freq_H number,
42                        freq_I number,
43                        freq_L number,
44                        freq_K number,
45                        freq_M number,
46                        freq_F number,
47                        freq_P number,
48                        freq_S number,
49                        freq_T number,
50                        freq_W number,
51                        freq_Y number,
52                        freq_V number
53                        )''')
54
55        cur.bindarraysize = 50
```

# The Original Honey Bee Protein Data Set

# The 1ˢᵗ CGI File (cont'd)

query_arraysize.sql | bind_insert.sql | Data_output_from_DB_dynamic.py | **Data_Input_dynamic.cgi** | notepatplus.html

```python
54
55          cur.bindarraysize = 50
56          cur.setinputsizes(10,9000,float,float,float,float,float,float,float,float,
57                            float,float,float,float,float,float,float,float,float,
58                            float,float, float)
59
60          #read raw data from a file
61          infile=file(theFile, 'r')
62          myStr=""
63          finalStr=''
64
65          #form a string with the raw data
66          for aline in infile:
67              myStr=myStr+aline
68
69          #form a continuous string
70          strL=myStr.replace('\n','')
71
72          #change the string into a list, one protein per list item
73          aList=strL.split('>')
74
75          #keep the list items that contains the substring, [Apis mellifera]
76          for anItem in aList:
77              if '[Apis mellifera]' in anItem:
78                  finalStr=finalStr+anItem
79
80          end=0
81          totalLength=len(finalStr)
```

# The Processed Bee Protein Data

# Slice out the Sub-strings of Interest from the String Containing the Bee Protein Data

gi|328788585|ref|XP_396734.4|

3

0    3                                12

PREDICTED: neurogenic locus Notch protein

[Apis mellifera]MNFFF...AIYIgi|328793854|

10

# The 1st CGI File (cont'd)

```python
80          end=0
81          totalLength=len(finalStr)
82          repetitions = finalStr.count('[Apis mellifera]')
83
84          #extract the target substrings, the gi number and the protein sequence
85          for i in range(repetitions):
86
87              start = finalStr.find('gi|', end) + 3
88              end = finalStr.find('|', start)
89              gi = finalStr[start : end]
90              #print 'gi = ', gi
91              start=finalStr.find('mellifera]', end)+10
92              end=finalStr.find('gi|', start)
93              if end == -1:
94                  end=totalLength
95              seq = finalStr[start:end]
96              #print 'seq =', seq
97              seqLength=len(seq)
98              freq_A=seq.count('A')/float(seqLength)
99              freq_R=seq.count('R')/float(seqLength)
100             freq_N=seq.count('N')/float(seqLength)
101             freq_D=seq.count('D')/float(seqLength)
102             freq_C=seq.count('C')/float(seqLength)
103             freq_Q=seq.count('Q')/float(seqLength)
104             freq_E=seq.count('E')/float(seqLength)
105             freq_G=seq.count('G')/float(seqLength)
106             freq_H=seq.count('H')/float(seqLength)
107             freq_I=seq.count('I')/float(seqLength)
```
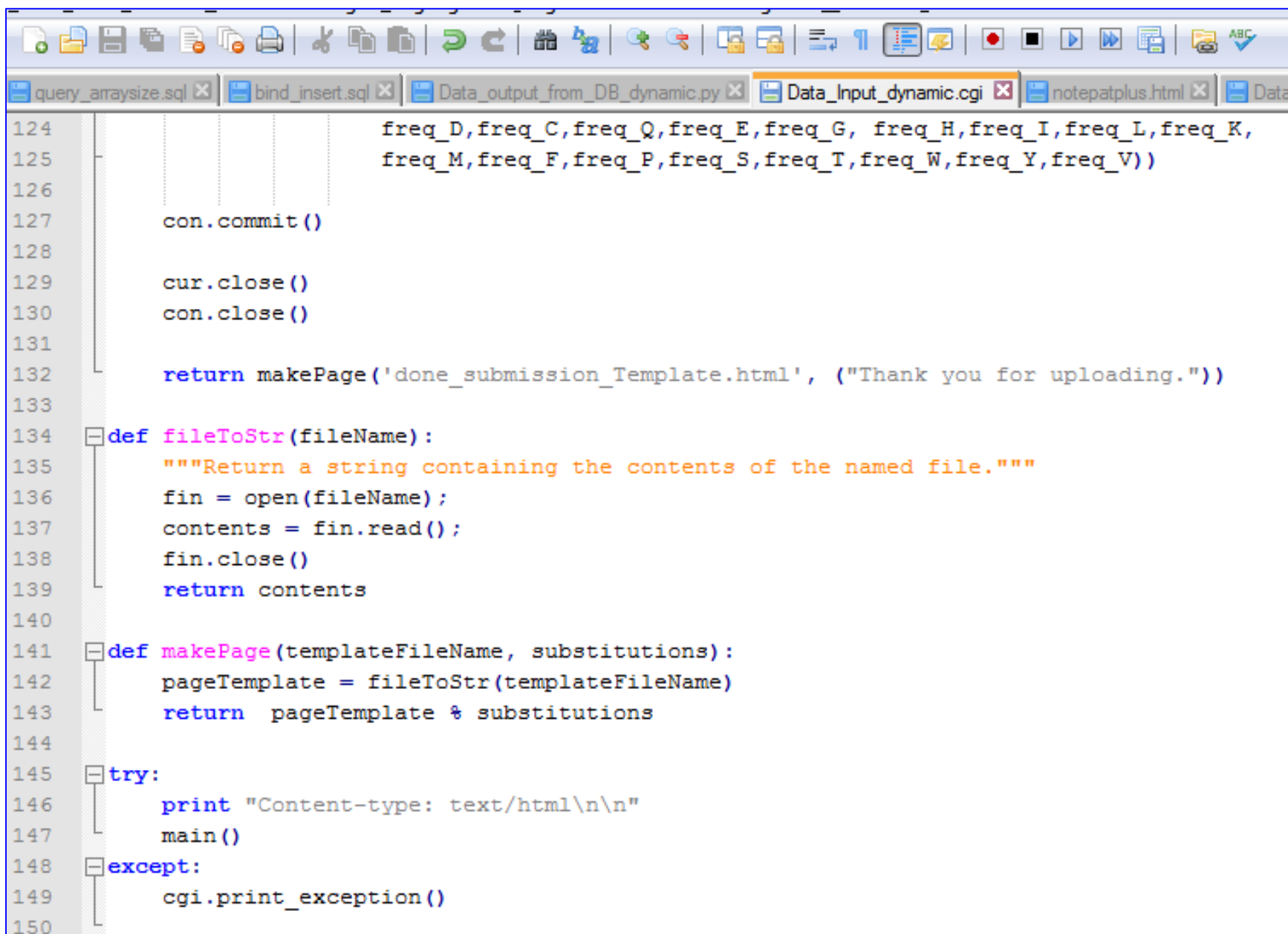
# The 1ˢᵗ CGI File (cont'd)

`query_arraysize.sql` | `bind_insert.sql` | `Data_output_from_DB_dynamic.py` | **Data_Input_dynamic.cgi** | `notepatplus.html` | `Data_out`

```python
106            freq_H=seq.count('H')/float(seqLength)
107            freq_I=seq.count('I')/float(seqLength)
108            freq_L=seq.count('L')/float(seqLength)
109            freq_K=seq.count('K')/float(seqLength)
110            freq_M=seq.count('M')/float(seqLength)
111            freq_F=seq.count('F')/float(seqLength)
112            freq_P=seq.count('P')/float(seqLength)
113            freq_S=seq.count('S')/float(seqLength)
114            freq_T=seq.count('T')/float(seqLength)
115            freq_W=seq.count('W')/float(seqLength)
116            freq_Y=seq.count('Y')/float(seqLength)
117            freq_V=seq.count('V')/float(seqLength)
118
119            cur.execute('''insert into beeProteins (gi, sequence, freq_A,freq_R,freq_N,
120                      freq_D,freq_C,freq_Q,freq_E,freq_G, freq_H,freq_I,freq_L,freq_K,
121                      freq_M,freq_F,freq_P,freq_S,freq_T,freq_W,freq_Y,freq_V) values(
122                      :v1,:v2,:v3,:v4,:v5,:v6,:v7,:v8,:v9,:v10,:v11,:v12,:v13,:v14,:v15,
123                      :v16,:v17,:v18,:v19,:v20,:v21,:v22)''',(gi,seq,freq_A,freq_R,freq_N,
124                      freq_D,freq_C,freq_Q,freq_E,freq_G, freq_H,freq_I,freq_L,freq_K,
125                      freq_M,freq_F,freq_P,freq_S,freq_T,freq_W,freq_Y,freq_V))
126
127        con.commit()
128
129        cur.close()
130        con.close()
131
132        return makePage('done_submission_Template.html', ("Thank you for uploading."))
```
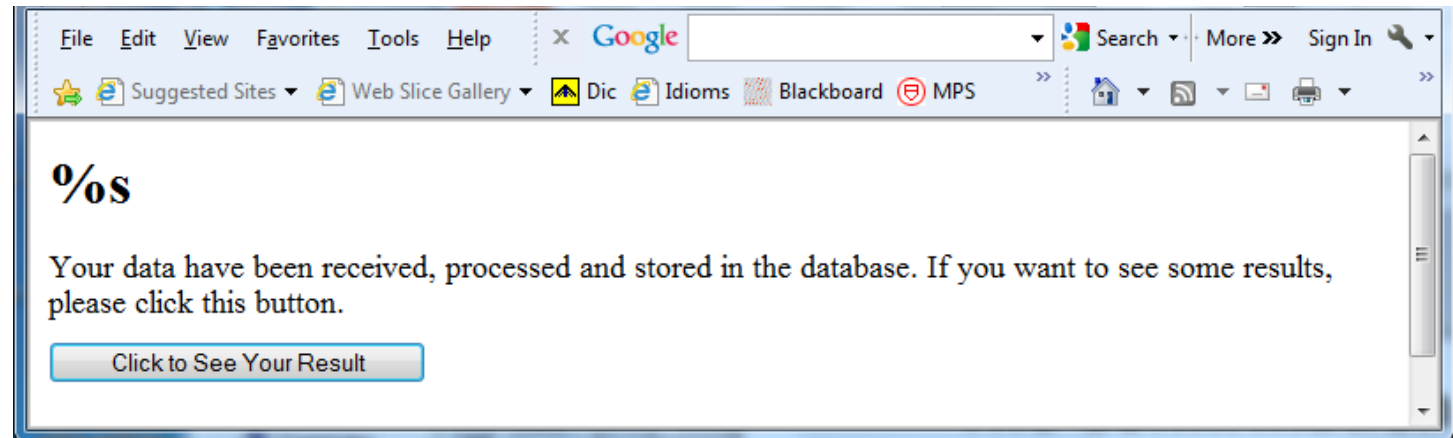
# The 1ˢᵗ CGI File (cont'd)

query_arraysize.sql | bind_insert.sql | Data_output_from_DB_dynamic.py | **Data_Input_dynamic.cgi** | notepatplus.html | Data
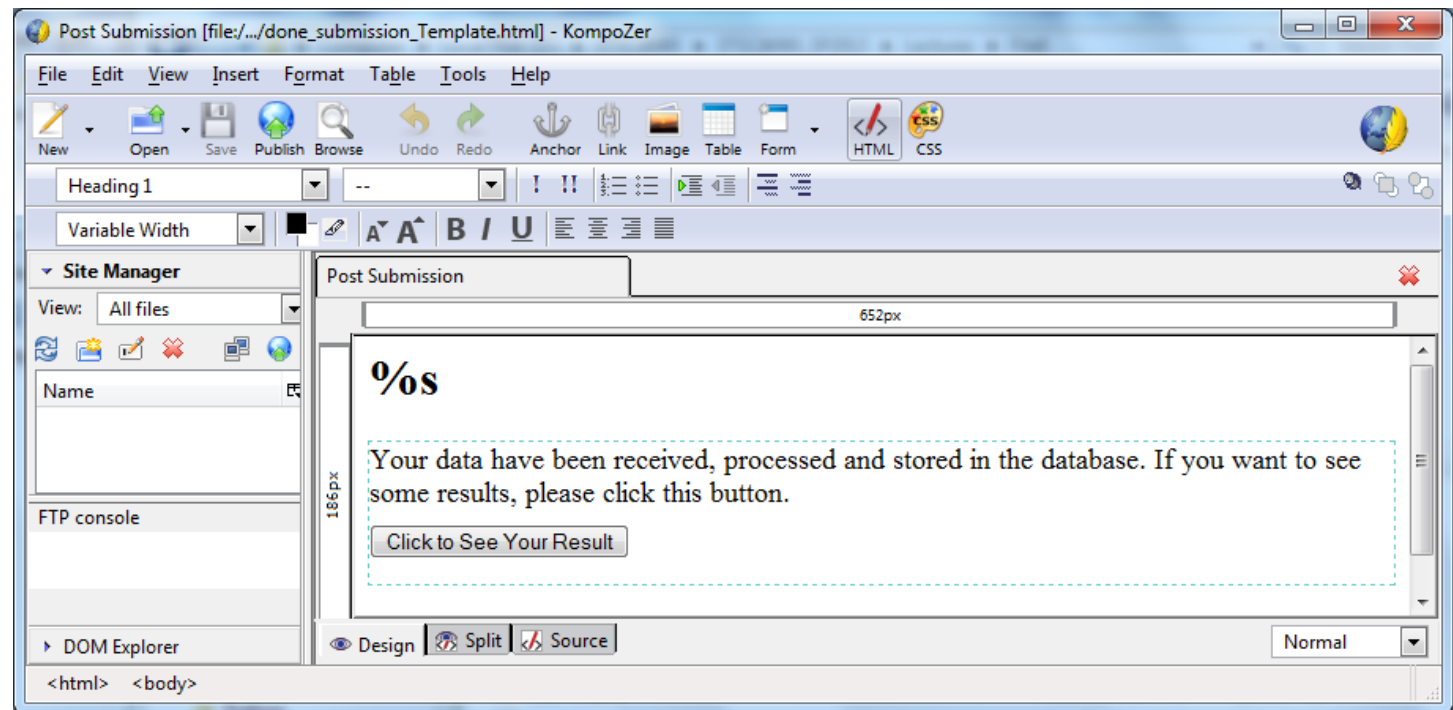
```
124                         freq_D,freq_C,freq_Q,freq_E,freq_G, freq_H,freq_I,freq_L,freq_K,
125                         freq_M,freq_F,freq_P,freq_S,freq_T,freq_W,freq_Y,freq_V))
126
127         con.commit()
128
129         cur.close()
130         con.close()
131
132         return makePage('done_submission_Template.html', ("Thank you for uploading."))
133
134 def fileToStr(fileName):
135     """Return a string containing the contents of the named file."""
136     fin = open(fileName);
137     contents = fin.read();
138     fin.close()
139     return contents
140
141 def makePage(templateFileName, substitutions):
142     pageTemplate = fileToStr(templateFileName)
143     return  pageTemplate % substitutions
144
145 try:
146     print "Content-type: text/html\n\n"
147     main()
148 except:
149     cgi.print_exception()
150
```

# The Confirmation Page Template

**IE View →**



**KompoZer Design View →**

# The 2$^{nd}$ CGI File

```
sec0.R    HW1-2    Assignment result    Data_output_from_DB_dynamic.cgi
 1      '''
 2      This program extracts the data from the Oracle database and
 3      display them on a web page upon user request.
 4
 5      ******** Coded by Xiaolong Yang @ Cornell University ********
 6      '''
 7
 8      import scipy as sp
 9      import cgi
10      import cx_Oracle
11
12      def main(): # NEW
13          #form=cgi.FieldStorage() #cgi script line
14          #theStr=form.getfirst('theList','')
15          contents = processInput()
16          print contents
17
18
19      def processInput(): #This function extracts data from a Oracle table.
20          con = cx_Oracle.connect('python/welcome')
21          cur=con.cursor()
22          aaList=['A','R','N','D','C','Q','E','G','H','I','L','K','M','F','P','S','T','W','Y','V']
23          fList=[() for t in range(20)]
24          for i in range(20):
25              myDict={'aa':aaList[i]}
26              obj=cur.execute('''select gi, freq_%(aa)s from beeproteins, (select max(freq_%(aa)s)
27                      as max%(aa)s from beeproteins) where freq_%(aa)s=max%(aa)s''' % myDict)
28              for x in obj:
```

# The 2<sup>nd</sup> CGI File (cont'd)

```
Data_output_from_DB_dynamic.py    Data_Input_dynamic.cgi    notepatplus.html    Data_output_from_DB_dynamic.cgi

28                    for x in obj:
29                        fList[i]=x
30
31          myTuple=()
32          for t in range(20):
33            myTuple = myTuple + fList[t]
34
35          cur.close()
36          con.close()
37
38          return makePage('see_result_template.html', myTuple)
39
40     def fileToStr(fileName):
41          """Return a string containing the contents of the named file."""
42          fin = open(fileName);
43          contents = fin.read();
44          fin.close()
45          return contents
46
47     def makePage(templateFileName, substitutions):
48          pageTemplate = fileToStr(templateFileName)
49          return  pageTemplate % substitutions
50
51     try:
52          print "Content-type: text/html\n\n"
53          main()
54     except:
55          cgi.print_exception()
```

# The Database-Driven Webpage Template: the Design View

see_result_template.html

**Site Manager**

View: All files

Name

Form Template

485px

## Honey Bee Proteins with Highest Occurance of Each Amino Acid

| Amino acid | Code | gi (Protein ID) | Highest Relative Frequency |
|---|---|---|---|
| Alanine | A | %s | %s |
| Arginine | R | %s | %s |
| Asparagine | N | %s | %s |
| Aspartic acid | D | %s | %s |
| Cysteine | C | %s | %s |
| Glutamic acid | E | %s | %s |
| Glutamine | Q | %s | %s |
| Glycine | G | %s | %s |
| Histidine | H | %s | %s |
| Isoleucine | I | %s | %s |
| Leucine | L | %s | %s |
| Lysine | K | %s | %s |
| Methionine | M | %s | %s |
| Phenylalanine | F | %s | %s |
| Proline | P | %s | %s |
| Serine | S | %s | %s |
| Threonine | T | %s | %s |
| Tryptophan | W | %s | %s |
| Tyrosine | Y | %s | %s |
| Valine | V | %s | %s |

634 px

FTP console

# The Database-Driven Webpage Template: the Source View

# The Database-Driven Webpage Template: the Source View (cont'd)

# User uploads a data file for processing

# Feedback: Successful data uploading and processing

# The Result of Database-driven, Python-control Webpage

**Honey Bee Proteins with Highest Occurance of Each Amino Acid**

| Amino acid | Code | gi (Protein ID) | Highest Relative Frequency |
|---|---|---|---|
| Alanine | A | 145386514 | 0.33606557377 |
| Arginine | R | 386649517 | 0.214788732394 |
| Asparagine | N | 58585170 | 0.137931034483 |
| Aspartic acid | D | 94158729 | 0.118421052632 |
| Cysteine | C | 328787226 | 0.128919860627 |
| Glutamic acid | E | 328783028 | 0.189090909091 |
| Glutamine | Q | 328792223 | 0.151660516605 |
| Glycine | G | 328781397 | 0.221991701245 |
| Histidine | H | 110755242 | 0.0696055684455 |
| Isoleucine | I | 94158668 | 0.177777777778 |
| Leucine | L | 110755367 | 0.183308494784 |
| Lysine | K | 328782196 | 0.160213618158 |
| Methionine | M | 58585126 | 0.0547945205479 |
| Phenylalanine | F | 328785536 | 0.100628930818 |
| Proline | P | 145386569 | 0.142857142857 |
| Serine | S | 386649511 | 0.19298245614 |
| Threonine | T | 213972576 | 0.1640625 |
| Tryptophan | W | 328785536 | 0.0440251572327 |
| Tyrosine | Y | 328781285 | 0.0706921944035 |
| Valine | V | 145386514 | 0.180327868852 |

Last update by Xiaolong Yang

# Analyze the Data Uploaded, Processed and Stored in the Oracle Database: Display a DB Table in SAS

```
libname beeprot oracle user='python'
password='welcome'
path = 'xe';
proc print data=beeprot.beeproteins;
run;
```

# Analyze the Data Uploaded, Processed and Stored in the Oracle Database: Cluster Analysis in SAS

```
libname beeprot oracle user='python'
password='welcome'
path = 'xe';

proc cluster data=beeprot.beeproteins
method=ward ccc pseudo
outtree=tree;
var freq_A freq_R freq_N freq_D freq_C
    freq_Q freq_E freq_G freq_H freq_I
    freq_L freq_K freq_M freq_F freq_P
    freq_S freq_T freq_W freq_Y freq_V;
run;
```

# Diagrams of CCC, PST2 and PSF

# Cluster Analysis Results: Values of CCC, PST2 and PSF

| Cluster History | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| Number of Clusters | Clusters Joined | | Freq | Semipartial R-Square | R-Square | Approximate Expected R-Square | Cubic Clustering Criterion | Pseudo F Statistic | Pseudo t-Squared | Tie |
| 461 | OB21 | OB25 | 2 | 0.0000 | 1.00 | . | . | . | . | T |
| 460 | OB93 | OB97 | 2 | 0.0000 | 1.00 | . | . | . | . | T |

| | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| 9 | CL12 | CL17 | 119 | 0.0195 | .487 | .440 | 8.63 | 53.7 | 14.0 | |
| 8 | CL48 | CL9 | 134 | 0.0229 | .464 | .426 | 6.47 | 56.1 | 15.4 | |
| 7 | CL10 | CL19 | 123 | 0.0232 | .441 | .409 | 5.46 | 59.7 | 18.1 | |
| 6 | CL8 | CL23 | 156 | 0.0250 | .416 | .389 | 4.72 | 64.8 | 15.5 | |
| **5** | CL11 | CL7 | 213 | 0.0532 | .362 | .363 | -.14 | **64.9** | **37.7** | |
| 4 | CL357 | CL26 | 85 | 0.0708 | .291 | .329 | -5.4 | 62.8 | 574 | |
| 3 | CL6 | CL4 | 241 | 0.0859 | .206 | .280 | -9.7 | 59.4 | 57.8 | |
| 2 | CL3 | CL86 | 249 | 0.0957 | .110 | .177 | -9.4 | 56.8 | 53.3 | |
| 1 | CL2 | CL5 | 462 | 0.1099 | .000 | .000 | 0.00 | . | 56.8 | |

# The Dendrogram