# Topics

› Review of basic data structures

› Accessing and working with objects in python

› Numpy

› Pandas
  – What are dataframes?
  – Sample column counting operation
  – Test, refine, test, refine, test, refine to get the data structure we want
  – Extracting a column

› Matplotlib
  – Making a simple line plot

› Scipy

# Next goal: Analyze base frequency at each position in the first 50 bases

› Need to iterate over sam file again

› Need to grab the first 50 bases of each read

› Need to turn the string into a list

› Need to move the data to a pandas dataframe
  – This is a more flexible structure for holding heterogenous data
  – Not as efficient as a numpy array, but still very fast
  – Very similar to dataframes in R language

› Need to count base occurrences in each column
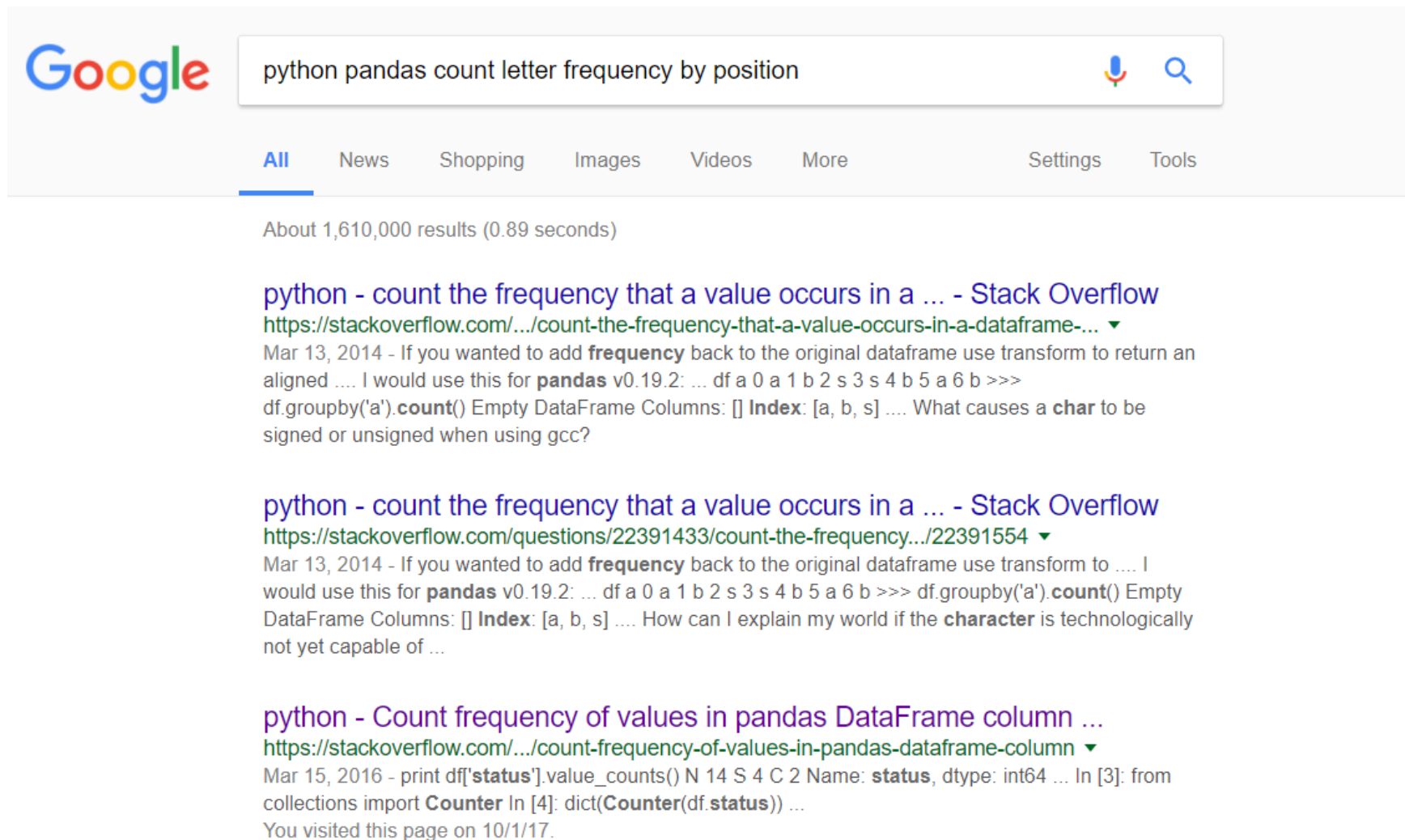
› New file: baseFrequencyCounter.py

# How to think of a dataframe?

```
In [23]: df4[(df4.Quantity == 1) & ((df4.State == 'NY') | (df4.State == 'CO'))]
```

Out[23]:

|    | ProductName | PaymentType | CustomerName | Quantity | UnitPrice | State |
|----|-------------|-------------|--------------|----------|-----------|-------|
| 0  | Product-A   | Visa        | Alice        | 1        | 44.99     | NY    |
| 5  | Product-B   | Visa        | Alice        | 1        | 14.99     | NY    |
| 6  | Product-F   | Discover    | John         | 1        | 29.49     | CO    |
| 13 | Product-H   | Visa        | Alice        | 1        | 79.99     | NY    |
| 14 | Product-A   | Discover    | John         | 1        | 40.49     | CO    |

# How to count letters in a column?

# Some test code

```python
readLengthMinimum = 100
analysisLength = 50

def getBaseDataFrame(samFile, readLengthMinimum, analysisLength, analysisPositionStart = 0, dataLimit = False):
    assert readLengthMinimum >= analysisLength + analysisPositionStart, "Analysis length must be shorter than or equal in length to the mini
    import samReader2
    import pandas
    samLines = samReader2.readSAMFileLines(samFile)
    baseTable = {}
    for line in samLines:
        if len(line.sequence.sequence) >= readLengthMinimum:
            if line.readID + ".1" in baseTable:  #dealing with paired end reads having duplicated indices
                baseTable[line.readID + ".2"] = list(line.sequence.sequence[analysisPositionStart:analysisLength + analysisPositionStart])
            else:
                baseTable[line.readID + ".1"] = list(line.sequence.sequence[analysisPositionStart:analysisLength + analysisPositionStart])
        if dataLimit:
            if len(qualityList) >= dataLimit:
                break
    baseDataFrame = pandas.DataFrame(baseTable)
    print(baseDataFrame)
    quit()
    return baseDataFrame

qualityMeanMatrix = getBaseDataFrame("sampleData.sam", 100, 50)
```

Data looks almost right, except we need it transposed (we want position to be the columns).

```
C:\Users\mweinstein\Documents\pythonClass2>python baseFrequencyCounter.py
Read 500172 lines
     SRR067577.10000004.1 SRR067577.10000004.2 SRR067577.1000004.1   \
0                       C                     C                     A
1                       C                     T                     A
2                       T                     T                     C
3                       A                     C                     C
4                       C                     T                     C
5                       C                     C                     A
6                       A                     T                     C
7                       G                     C                     T
8                       A                     T                     T
9                       C                     G                     C
10                      C                     C                     T
11                      G                     A                     A
12                      G                     G                     C
13                      C                     A                     C
14                      T                     A                     C
36                      T                     T                     A
37                      T                     C                     T
38                      A                     A                     C
39                      C                     C                     A
40                      A                     C                     G
41                      G                     A                     C
42                      T                     G                     A
43                      T                     C                     G
44                      A                     C                     G
45                      G                     A                     A
46                      A                     T                     T
47                      G                     G                     G
48                      A                     T                     A
49                      A                     A                     G

[50 rows x 500000 columns]
```
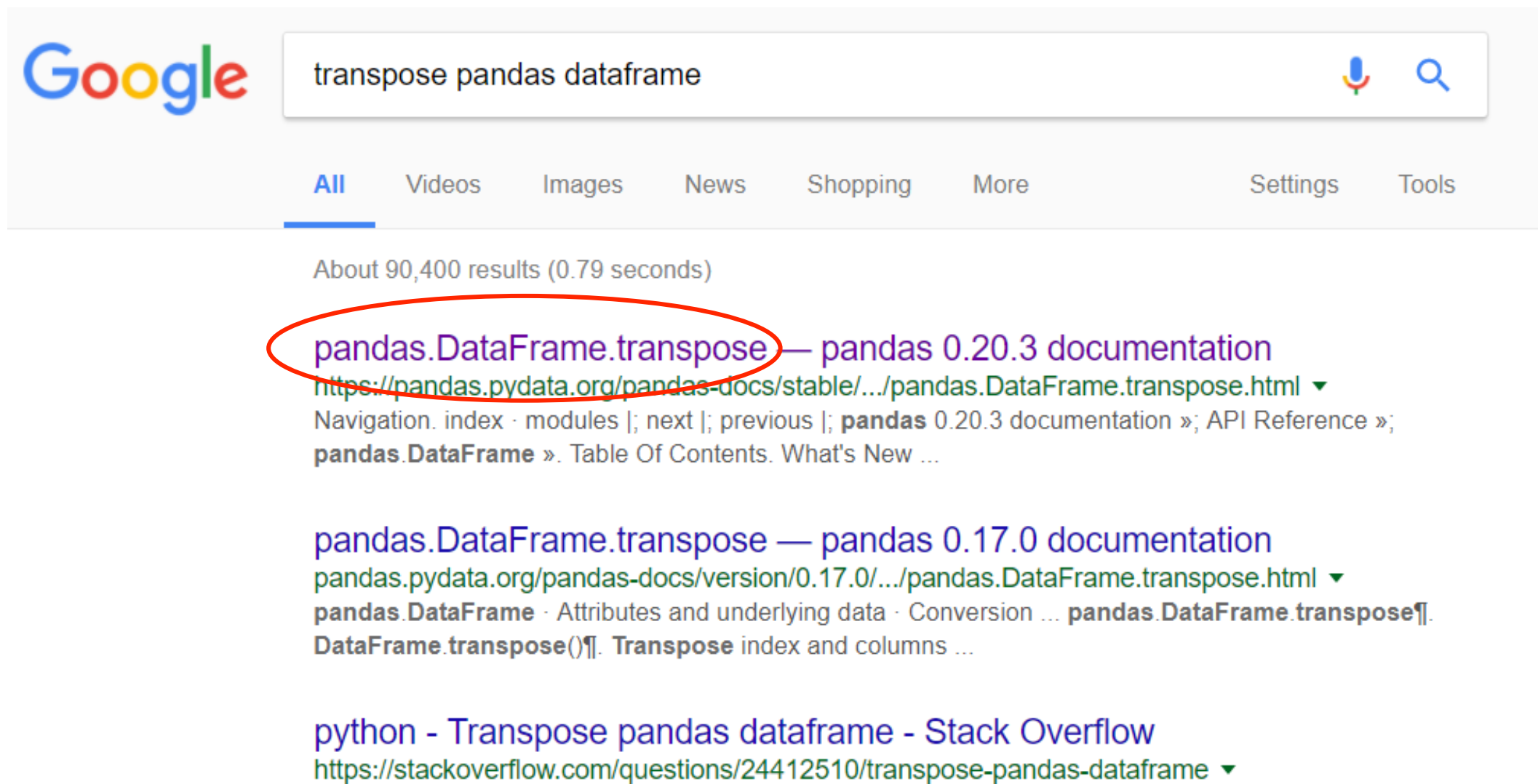
# How to transpose

# Testing our data structure

```python
readLengthMinimum = 100
analysisLength = 50

def getBaseDataFrame(samFile, readLengthMinimum, analysisLength, analysisPositionStart = 0, dataLimit = False):
    assert readLengthMinimum >= analysisLength + analysisPositionStart, "Analysis length must be shorter than or equal in length to the min
    import samReader2
    import pandas
    samLines = samReader2.readSAMFileLines(samFile)
    baseTable = {}
    for line in samLines:
        if len(line.sequence.sequence) >= readLengthMinimum:
            if line.readID + ".1" in baseTable:  #dealing with paired end reads having duplicated indices
                baseTable[line.readID + ".2"] = list(line.sequence.sequence[analysisPositionStart:analysisLength + analysisPositionStart])
            else:
                baseTable[line.readID + ".1"] = list(line.sequence.sequence[analysisPositionStart:analysisLength + analysisPositionStart])
        if dataLimit:
            if len(qualityList) >= dataLimit:
                break
    baseDataFrame = pandas.DataFrame(baseTable)
    print(baseDataFrame)
    baseDataFrame = baseDataFrame.transpose()
    print(baseDataFrame)
    test = baseDataFrame[0].value_counts()
    print(test)
    test = dict(test)
    print(test)
    quit()
    return baseDataFrame

qualityMeanMatrix = getBaseDataFrame("sampleData.sam", 100, 50)
```

```
[50 rows x 500000 columns]
                      0  1  2  3  4  5  6  7  8  9  ... 40 41 42 43 44 45 46  \
SRR067577.10000004.1  C  C  T  A  C  C  A  G  A  C ...  C  T  C  C  C  A  A
SRR067577.10000004.2  C  T  T  C  T  C  T  C  T  G ...  A  C  A  A  C  C  A
SRR067577.1000004.1   A  A  C  C  C  A  C  T  T  C ...  A  T  G  T  A  A  C
SRR067577.1000004.2   A  T  A  T  A  T  A  C  A  T ...  C  T  C  C  C  T  C
SRR067577.10000049.1  T  T  A  T  T  A  A  T  G  A ...  T  C  T  T  A  C  T
SRR067577.10000049.2  T  A  C  T  T  C  T  T  G  T ...  C  A  C  A  T  G  G
SRR067577.10000090.1  A  G  C  C  T  A  C  G  A  G ...  C  A  G  G  A  C  T
SRR067577.10000090.2  G  C  T  C  T  G  G  G  C  A ...  C  A  A  A  A  T  A
SRR067577.10000166.1  A  G  G  G  A  A  G  G  C  T ...  A  A  G  A  G  G  G
SRR067577.10000166.2  G  C  A  G  T  C  G  A  C  T ...  C  C  T  G  T  A  A
SRR067577.1000018.1   G  C  C  C  A  C  C  T  T  G ...  A  G  C  C  A  C  C
SRR067577.1000018.2   C  C  C  G  C  C  T  T  G  G ...  G  C  C  A  C  C  A
SRR067577.10000199.1  C  A  T  G  A  G  C  T  G  G ...  G  A  G  C  T  G  G
SRR067577.10000199.2  G  T  G  A  G  A  A  G  A  G ...  C  A  G  G  C  A  A
SRR067577.10000207.1  A  G  A  A  T  G  A  T  G  A ...  C  A  G  T  C  T  A
SRR067577.10000207.2  T  C  A  T  A  A  C  A  G  C ...  A  G  A  T  G  T  C
SRR067577.10000211.1  C  T  T  G  T  T  T  T  C  T ...  A  A  A  T  G  A  A
SRR067577.10000211.2  T  A  A  G  T  T  C  A  T  T ...  A  G  A  G  C  C  T
SRR067577.10000216.1  T  C  C  T  C  G  T  C  C  T ...  T  C  C  C  C  G  T
SRR067577.10000216.2  C  C  C  A  T  C  G  C  C  A ...  C  C  C  C  A  G  T
SRR067577.10000224.1  C  T  C  C  C  A  G  G  A  C ...  A  C  A  G  A  G  C
SRR067577.10000224.2  T  G  C  T  T  T  C  A  T  G ...  T  T  T  C  T  G  C
SRR067577.10000326.1  A  T  T  G  T  A  T  A  C  A ...  C  C  A  C  C  A  A
SRR067577.10000326.2  G  A  A  G  A  G  A  C  A  G ...  C  C  C  T  G  G  G
SRR067577.10000377.1  C  A  T  T  C  T  T  A  C  A ...  C  A  C  T  G  A  T
SRR067577.10000377.2  T  A  T  C  C  A  T  G  T  T ...  G  T  T  T  T  C  T
SRR067577.10000428.1  G  T  T  C  A  G  G  C  C  A ...  A  A  G  A  G  G  T
SRR067577.10000428.2  A  C  T  C  T  G  A  A  A  T ...  G  G  T  G  A  G  A
SRR067577.10000476.1  C  A  A  C  T  G  C  A  A  C ...  A  A  A  G  A  G  C
SRR067577.10000476.2  A  G  C  A  T  T  A  T  T  G ...  G  G  A  A  T  A  G
...                   .. .. .. .. .. .. .. .. .. ..      .. .. .. .. .. .. ..
SRR067577.9999566.1   A  T  T  T  T  T  T  A  A  A ...  T  T  G  T  A  G  T
SRR067577.9999566.2   A  G  T  T  T  T  A  T  T  C ...  G  A  A  T  G  T  A
SRR067577.9999624.1   G  G  G  C  T  A  G  G  T  C ...  A  A  G  G  G  T  C
```

```
                      47 48 49
SRR067577.10000004.1   C  G  A
SRR067577.10000004.2   A  T  T
SRR067577.1000004.1    C  A  G
SRR067577.1000004.2    C  A  T
SRR067577.10000049.1   T  T  C
SRR067577.10000049.2   T  T  A
SRR067577.10000090.1   G  C  C
SRR067577.10000090.2   T  G  A
SRR067577.10000166.1   A  G  A
SRR067577.10000166.2   T  G  C
SRR067577.1000018.1    A  T  G
SRR067577.1000018.2    T  G  C
SRR067577.10000199.1   A  A  T
SRR067577.10000199.2   G  A  G
SRR067577.10000207.1   A  G  G
SRR067577.10000207.2   A  A  G
SRR067577.10000211.1   T  T  T
SRR067577.10000211.2   C  A  C
SRR067577.10000216.1   G  G  C
SRR067577.10000216.2   G  C  T
SRR067577.10000224.1   T  G  T
SRR067577.10000224.2   A  A  A
SRR067577.10000326.1   T  G  C
SRR067577.10000326.2   A  T  T
SRR067577.10000377.1   T  T  G
SRR067577.10000377.2   C  C  T
SRR067577.10000428.1   T  T  G
SRR067577.10000428.2   G  A  G
SRR067577.10000476.1   A  T  T
SRR067577.10000476.2   A  T  G
...                    .. .. ..
SRR067577.9999566.1    G  G  A
SRR067577.9999566.2    G  T  A
SRR067577.9999624.1    A  C  A
```

# Looks good!

```
[500000 rows x 50 columns]
T      147458
A      140610
C      107953
G      103899
N          80
Name: 0, dtype: int64
{'T': 147458, 'A': 140610, 'C': 107953, 'G': 103899, 'N': 80}
```

# Testing our data structure

```
3    readLengthMinimum = 100
4    analysisLength = 50
5
6    def getBaseDataFrame(samFile, readLengthMinimum, analysisLength, analysisPositionStart = 0, dataLimit = False):
7        assert readLengthMinimum >= analysisLength + analysisPositionStart, "Analysis length must be shorter than or equal in length to the min
8        import samReader2
9        import pandas
10       samLines = samReader2.readSAMFileLines(samFile)
11       baseTable = {}
12       for line in samLines:
13           if len(line.sequence.sequence) >= readLengthMinimum:
14               if line.readID + ".1" in baseTable:  #dealing with paired end reads having duplicated indices
15                   baseTable[line.readID + ".2"] = list(line.sequence.sequence[analysisPositionStart:analysisLength + analysisPositionStart])
16               else:
17                   baseTable[line.readID + ".1"] = list(line.sequence.sequence[analysisPositionStart:analysisLength + analysisPositionStart])
18           if dataLimit:
19               if len(qualityList) >= dataLimit:
20                   break
21       baseDataFrame = pandas.DataFrame(baseTable)
22       print(baseDataFrame)
23       baseDataFrame = baseDataFrame.transpose()
24       print(baseDataFrame)
25       test = baseDataFrame[0].value_counts()
26       print(test)
27       test = dict(test)
28       print(test)
29       quit()
30       return baseDataFrame
31
32   qualityMeanMatrix = getBaseDataFrame("sampleData.sam", 100, 50)
```

<span style="color:red">Test code (can be removed after it works)<br>We just need to return the transposed dataframe<br>Replace the box contents with:<br>return baseDataFrame.transpose()</span>

# Testing out our next data structure

```python
23
24  def extractBaseCountDataFrame(baseDataFrame):
25      import pandas
26      indices = list(baseDataFrame.columns.values)
27      baseCountTable = {}
28      for index in indices:
29          print(index)
30          baseCountTable[index] = baseDataFrame[index].value_counts()
31      baseCountDataFrame = pandas.DataFrame(baseCountTable)
32      print(baseCountDataFrame)
33      quit()
34
35  baseDataFrame = getBaseDataFrame("sampleData.sam", 100, 50)
36  extractBaseCountDataFrame(baseDataFrame)
37
```

# Testing out the data structure

# Transposing and extracting a column...

```python
def extractBaseCountDataFrame(baseDataFrame):
    import pandas
    indices = list(baseDataFrame.columns.values)
    baseCountTable = {}
    for index in indices:
        print(index)
        baseCountTable[index] = baseDataFrame[index].value_counts()
    baseCountDataFrame = pandas.DataFrame(baseCountTable)
    print(baseCountDataFrame)
    baseCountDataFrame = baseCountDataFrame.transpose()
    print(baseCountDataFrame["A"])
    quit()


baseDataFrame = getBaseDataFrame("sampleData.sam", 100, 50)
extractBaseCountDataFrame(baseDataFrame)
```

This is what we want to see

```
[5 rows x 50 columns]
0      140610.0
1      145615.0
2      143508.0
3      142440.0
4      143105.0
5      143072.0
6      143703.0
7      143948.0
8      143004.0
9      143183.0
10     143166.0
11     143227.0
12     143630.0
13     143039.0
14     143337.0
15     143288.0
16     143769.0
17     143421.0
18     143207.0
19     143117.0
20     143231.0
21     142820.0
22     143374.0
23     142784.0
24     143335.0
25     143396.0
26     143561.0
27     143462.0
28     143308.0
29     142722.0
30     143586.0
31     143822.0
32     143662.0
```

# Finalize the function and call

```python
24  def extractBaseCountDataFrame(baseDataFrame):
25      import pandas
26      indices = list(baseDataFrame.columns.values)
27      baseCountTable = {}
28      for index in indices:
29          baseCountTable[index] = baseDataFrame[index].value_counts()
30      baseCountDataFrame = pandas.DataFrame(baseCountTable)
31      return baseCountDataFrame.transpose()
32
33  baseDataFrame = getBaseDataFrame("sampleData.sam", readLengthMinimum, analysisLength)
34  baseCountDataFrame = extractBaseCountDataFrame(baseDataFrame)
35
```

# Important question:
# How to represent our data

› Class participation time: Start making suggestions

› Don't cheat here if you've looked ahead

› More than one right answer

# Pyplot marker codes
## (Don't worry about memorizing)
## Keep a reference if you use these often

| character | color |
|-----------|-------|
| 'b' | blue |
| 'g' | green |
| 'r' | red |
| 'c' | cyan |
| 'm' | magenta |
| 'y' | yellow |
| 'k' | black |
| 'w' | white |

| character | description |
|-----------|-------------|
| '-' | solid line style |
| '--' | dashed line style |
| '-.' | dash-dot line style |
| ':' | dotted line style |
| '.' | point marker |
| ',' | pixel marker |
| 'o' | circle marker |
| 'v' | triangle_down marker |
| '^' | triangle_up marker |
| '<' | triangle_left marker |
| '>' | triangle_right marker |
| '1' | tri_down marker |
| '2' | tri_up marker |
| '3' | tri_left marker |
| '4' | tri_right marker |
| 's' | square marker |
| 'p' | pentagon marker |
| '*' | star marker |
| 'h' | hexagon1 marker |
| 'H' | hexagon2 marker |
| '+' | plus marker |
| 'x' | x marker |
| 'D' | diamond marker |
| 'd' | thin_diamond marker |
| '|' | vline marker |
| '_' | hline marker |

# Iterative plotting… it's really this simple

```
32
33  def plotBaseCounts(baseCountDataFrame, length):
34      import matplotlib.pyplot as plt
35      baseMarker = {"A" : "b-",
36                    "C" : "r-",
37                    "G" : "g-",
38                    "T" : "y-",
39                    "N" : "k-"}
40      for base in "ATGCN":
41          plt.plot(range(length), baseCountDataFrame[base], baseMarker[base], label = base)
42      plt.xlabel("Position")
43      plt.ylabel("Counts")
44      plt.legend()
45      #plt.show()
46      plt.savefig("baseCounts.png")
47
48  baseDataFrame = getBaseDataFrame("sampleData.sam", readLengthMinimum, analysisLength)
49  baseCountDataFrame = extractBaseCountDataFrame(baseDataFrame)
50  plotBaseCounts(baseCountDataFrame, analysisLength)
```

Success!