

$$Q1 \Rightarrow p = 0.25 \quad 1-p = 0.75$$

$$n = 1000$$

$$n = 240$$

Central Limit Theorem:

$$\phi(z) = \phi\left(\frac{[n - \mu]}{\sigma}\right) = \phi\left(\frac{[240 - \mu]}{\sigma}\right)$$

$$\sigma = \sqrt{np(1-p)} = 13.69$$

$$\mu = \text{mean} = n \cdot p = E[x] = 250$$

So

$$z = \left| \frac{240 - 250}{13.69} \right| \approx -0.7303$$

$$\text{So } \phi(z) = \phi(-0.7303) = \underline{0.7673}$$

Manual Approximate calculation = 0.7673

Simulate Probability = 0.7767

Calculated using scipy = 0.7776521016

codes are provides as screenhots

[28] ▶ MI

```
# Question 1
import numpy as np
import matplotlib.pyplot as plt
from scipy import stats
import random
import math

p = 0.25
n = 1000

BinomScipy = stats.binom(n,p)

random.seed(0)

arr = ['A','T','C','G'] # Since p = 0.25 , here also p(A) = 0.25

x = []
y = []

temp = 0
temp_arr = []
temp_1 = 0
for j in range(0,10000):
    temp = 0
    for k in range(0,1000):
        temp += 1 if (random.choice(arr) == 'A') else 0
    if(temp >= 240):
        temp_1 += 1

print("Simulated Probability of A appearing atleast 240 time : " + str(temp_1/10000) )
print("Calculated Probability of A appearing atleast 240 time : " + str(1-BinomScipy.cdf(239.9999)))
```

Simulated Probability of A appearing atleast 240 time : 0.7767  
Calculated Probability of A appearing atleast 240 time : 0.7776521016977749

$$\underline{Q \# 2} \Rightarrow p = 0.3 \quad n = 10$$

$$P(X=k) = {}^n C_k \cdot p^k \cdot (1-p)^{n-k}$$

$$P(X=0) = {}^{10} C_0 \cdot (0.3)^0 \cdot (1-0.3)^{10-0}$$

$$= \underline{0.0282}$$

$$P(X=2) = {}^{10} C_2 \cdot (0.3)^2 \cdot (1-0.3)^8$$

$$= \underline{0.2334}$$

$$E(X) = n \cdot p = 10 \times 0.3 = \underline{3}$$

$$\text{Var}[X] = n \cdot p \cdot (1-p) = 10 \times 0.3 \times (1-0.3) \\ = \underline{\underline{2.099}}$$



Q #3  $\Rightarrow$  a) k-mer based approach is used to check the similarities and dissimilarities in different biological sequences.

b) k-mer has extensive use in biotechnology in controlling translational efficiency.

c) They are also used to create attenuated vaccines.

d) They also have application in identifying heterogeneous gene expression.

e) They are used in construction of De Bruijn graphs.

f)  $k$ -mer distributions are well-preserved,  
Because of this reason bacterial genomes  
can be clustered into natural groups  
according to  $k$ -mer distribution similarities.

g) observed frequencies of  $k$ -words can be  
used to make inferences about DNA  
sequences.



Q# 4  $\Rightarrow$

	G	G	C	T	G	C	A	A	C	T	A	G	C	T	C
G	x	x			x							x			
G	x	x			x							x			
G	x	x			x							x			
T				x						x				x	
A							x	x			x				
A							x	x			x				
G	x	x			x							x			
C			x			x			x				x		x
T				x						x				x	
T				x						x				x	
G	x	x			x							x			
C			x			x			x				x		x

Conserved Regions: The region of similarity between 2 genomes

We use dotplot to identify conserved regions that are showing a match by x wherever a base matches or a group of bases matches

1.ipynb M

question5.py X

home > bhaeyan > question5.py > ...

```
1  # Question #5 Code
2  import numpy
3  import matplotlib.pyplot as plt
4
5  def df(x,y):
6      return 0 if x == y else 1
7
8  def M(seq1,seq2,i,j,k):
9      return sum(df(x,y) for x,y in zip(seq1[i:i+k],seq2[j:j+k]))
10
11 def ML(seq1,seq2,k):
12     n = len(seq1)
13     m = len(seq2)
14     return [[M(seq1,seq2,i,j,k) for j in range(m-k+1)] for i in range(n-k+1)]
15
16 def drawMatrix(M,t, seq1, seq2, nonblank = chr(0x25A0), blank = ' |'):
17     print(' ',end=" ")
18     for i in seq2:
19         print("|" + i,end=" ")
20     print("")
21     print('-'*(2 + 2*len(seq2)))
22     for label,row in zip(seq1,M):
23         line = ''.join('X|' if s < t else blank for s in row)
24         print(label + '|' + line)
25         # print('-'*(2 + 2*len(seq2)))
26
27 def dotplot(seq1,seq2,k = 1,t = 1):
28     M = ML(seq1,seq2,k)
29     drawMatrix(M, t, seq1,seq2)
30
31 # Sequence :
32 seqy = seqx = 'TGGCACACTCACACCACACAGACAGTTA'
33
34 # Plotting Function
35 dotplot(seqx,seqy)
```



```
~> python question5.py
```

```
|T|G|G|C|A|C|A|C|T|C|A|C|A|C|C|A|C|A|C|A|G|A|C|A|G|T|T|A
```

T	X						X																	X	X
G		X	X																X				X		
G		X	X															X				X			
C				X		X	X	X	X	X	X	X	X	X	X	X	X			X					
A					X			X	X	X		X	X	X	X	X	X	X	X	X	X			X	
C					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
A					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
C	X				X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
T	X						X															X	X		
C				X		X	X	X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
A					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
C					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
A					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
C					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
A					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
C					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
A					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
C					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
A					X			X	X	X	X	X	X	X	X	X	X	X	X	X	X			X	
G		X	X															X				X			
A				X		X												X				X			
C				X		X												X				X			
A				X		X												X				X			
G		X	X															X				X			
T	X						X															X	X		
T	X						X															X	X		
A				X		X			X	X		X	X	X	X	X	X	X	X	X	X			X	

```
~>
```





Diagonals representing complementary regions are shown from top-right left and bottom right diagonal directions.



RNA & its reverse used

Palindrome sequence

i) AU

ii) GC

iii) UGG CAUGCCA

AUGUGGCAUGCCA

RNA  $\Rightarrow$  AUGUGGCAUGCCA

Reverse  $\Rightarrow$  CCUGGCAUGCCA  
complement