

▼ Lab 1

```
1  # importing
2  import pandas as pd
3  import numpy as np
4  import matplotlib.pyplot as plt
5  import seaborn as sns
6  import math as m
7
8  # Bernoulli
9  def bernoulli(success_event, prob):
10     return (prob**success_event)*((1-prob)**(1-success_event))
11
12 def plot_bernoulli():
13     # successful and unsuccessful event
14     x = [0,1]
15     prob = 0.2
16     y = [None]*2
17     for _ in range(len(x)):
18         y[_] = bernoulli(x[_],prob)
19
20     plt.bar(x,y)
21
22 def normal(x,mean, std):
23     return (1/std*(np.pi)**(1/2)) * ((np.e)**(-(x-mean)**2/2*(std**2)))
24
25 def plot_normal():
26     x = [-1,-1,0,0,0,0,0,0.5,1,1]
27
28     m, std = 0,1
29     y = [None]*len(x)
30     for i in range(len(x)):
31         y[i] = normal(x[i],m,std)
32
33     normal_dist = pd.DataFrame({'x':x,'y':y})
34
35     plt.bar(x,y)
36
37 def factorial(n):
38     prod = 1
39     for i in range(2,n+1):
40         prod *= i
41     return prod
42
43 def binomial(x,p,q,n):
44
45
46
47     c_term = factorial(n) / (factorial(x) * factorial((n-x)))
48
49     return c_term*(p**x)*((1-p)**(1-x))
50
51
52
53 def plot_binomial():
54
55     x = [i for i in range(100)]
56     n = 500
57     p = 0.2
58     q = 1 - p
59
60     y = [None]*len(x)
61
62     for i in range(len(x)):
63         y[i] = binomial(x[i],p,q,n)
64     plt.bar(x,y)
65
66 def poisson(x,lamb):
67     return (m.e)**(-(lamb))*(lamb**x)/factorial(x)
68
69 def plot_poisson():
70     x = [0,1,2,3,4,5]
71     r = 0.5
72     y = [None]*len(x)
73
74     for i in range(len(x)):
75         y[i] = poisson(x[i],r)
76     plt.bar(x,y)
```

```

76     plt.bar(x,y)
77
78     def expo(x,lamb):
79         if x<=0:
80             return 0
81
82         return lamb*((m.e)**(-(lamb*x)))
83
84     def plot_expo():
85
86         x = np.linspace(-1,5,60)
87         x = x.tolist()
88         r = 0.5
89         y = [None]*len(x)
90
91         for i in range(len(x)):
92             y[i] = expo(x[i],r)
93         plt.bar(x,y)
94
95

```

▼ Lab 2

1. alpha: probability that fail to reject H0 when true
2. p-value - probability that get extreme values than calculated z-stat assuming H0 True (randomness)
3. if $p < \alpha$, reject otherwise failed to reject : more randomness $p = 0.05$: 5 % chance that Null is True , doesn't mean 95% Ha is correct
4. One Tailed , Two tailed, one sample , two sample
5. 5 % alpha, 95% confidence - z is ± 1.96 for two tailed
6. $z = \frac{x - u}{\text{std}/\sqrt{n}}$ for 2 sample - $\frac{s_1 - s_2}{\sqrt{\frac{\text{std}_1^2}{n_1} + \frac{\text{std}_2^2}{n_2}}}$

▼ Lab 3

```

1 filename = 'DATA.txt'
2 with open(filename,'r') as f:
3     lines = f.read().replace("\n", " ")
4     objs = lines.split()
5

```

```
'A B C A B C B D A C D D D C C '
```

General Points

BDA Domain Process

1. Databases , DFS, Paralle Programming
2. ML
3. Cloud
4. Apache Hive - Parallel Programming , HiveQL Scalaable, natural queries into mapreduce jobs over clusters - distributed computing, metadata indexes - on top of HDFS - METADATA - ACCESS CONTROL, PARTITION EFFICIENCY ,PROCESSING SPEED
5. Data science - exact and analytical, BDA is not simple , requires optimization , scalability, large datasets , no exact method right, combination of methods

Could not connect to the reCAPTCHA service. Please check your internet connection and reload to get a reCAPTCHA challenge.

