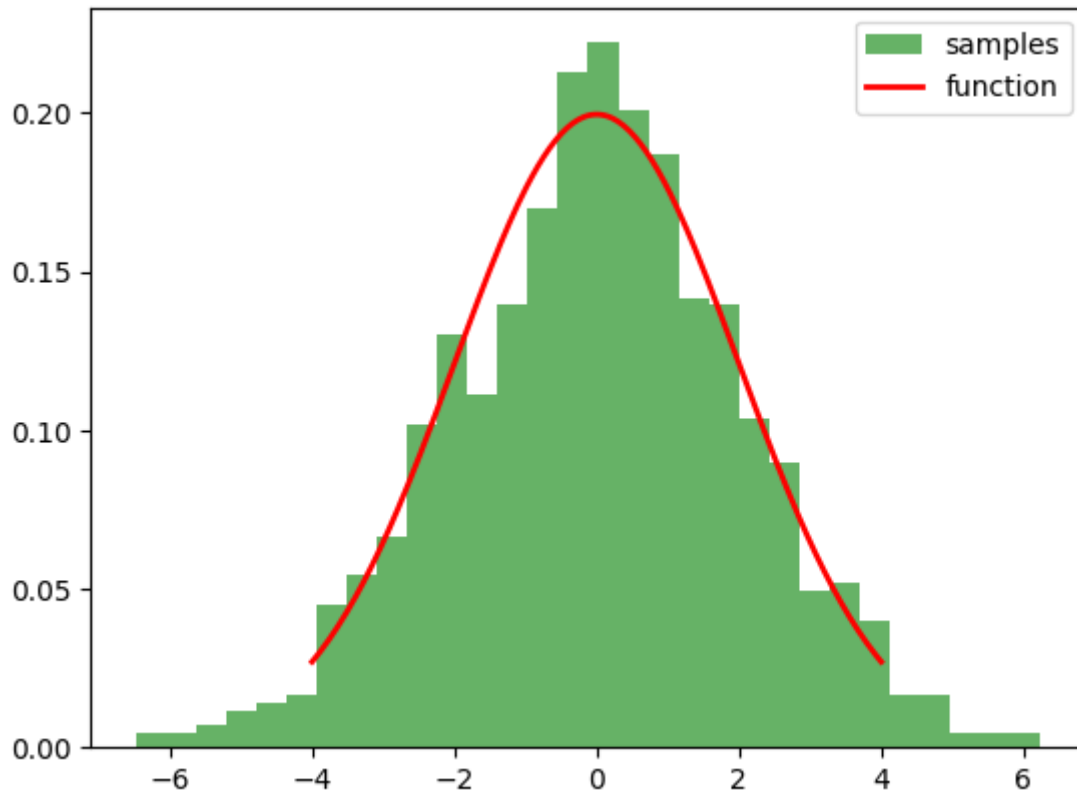


1.正态分布

随机变量 X 符合正态分布记为: $X \sim N(\mu, \sigma^2)$ 其中 $E(X) = \mu$, 方差为 σ^2

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 import scipy as sp
4
5 # defin c mu sigma
6 mu = 0
7 sigma = 2
8
9 # 随机分布样本 numpy.random.normal(mu,sigma,nums of points)
10 samples = np.random.normal(mu,sigma,1000)
11
12 plt.hist(samples, bins=30, density=True, alpha=0.6, color='g',
13 label='samples')
14
15 # 概率密度函数 scipy.stats.norm.pdf(X,mu,sigma)
16 x = np.linspace(-4,4,200)
17 pdf = sp.stats.norm.pdf(x,mu,sigma)
18 plt.plot(x, pdf, 'r', lw=2, label='function')
19 plt.legend()
20 plt.show()
```



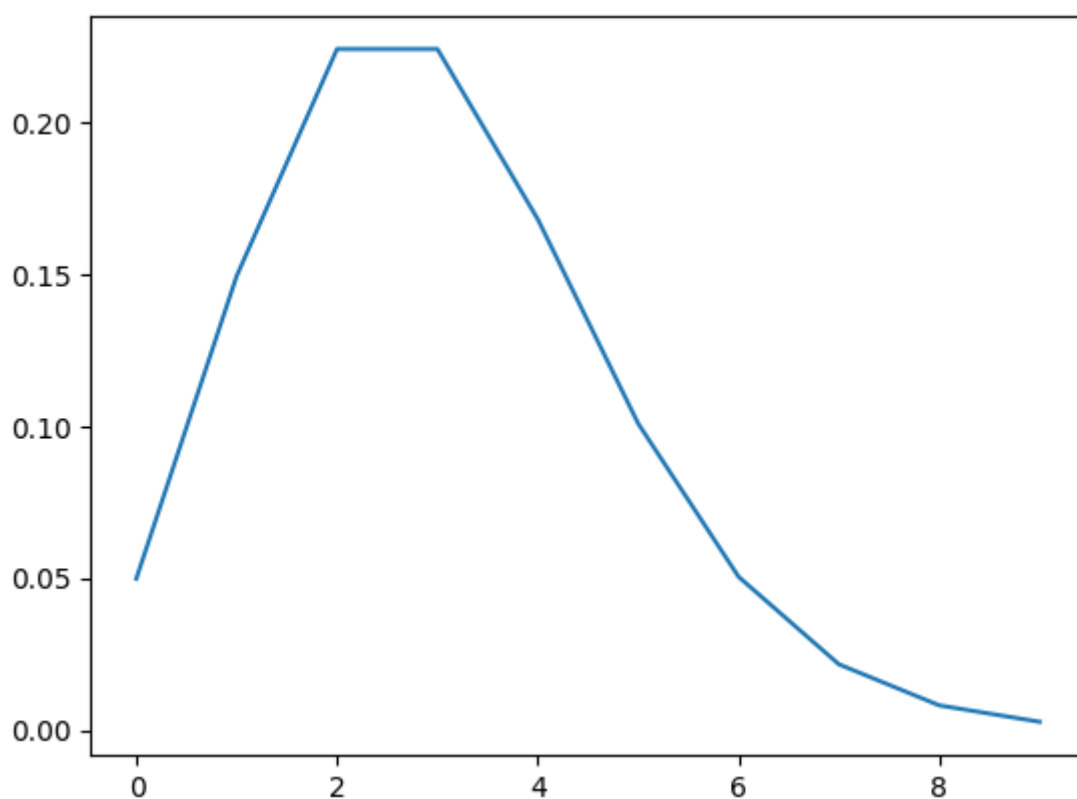
2.泊松分布 poisson

随机变量 X 符合参数为 λ 的泊松分布记为: $X \sim P(\lambda)$,Poisson分布的数学期望和方差都为 λ .其概率密度函数为

$$P(X = k) = \frac{\lambda^k}{k!} e^{-\lambda}, k = 1, 2, \dots, \lambda > 0$$

```
1 # poisson分布随机数生成:scipy.stats.poisson.rvs(lambda,size)
2 lambda_poisson = 3
3 X = sp.stats.poisson.rvs(lambda_poisson, size = 20)
4 print(X)
5
6 # poisson分布概率密度函数:scipy.stats.poisson.pmf(X,lambda)
7 k = np.arange(0,10)
8 probability = sp.stats.poisson.pmf(k, lambda_poisson)
9 plt.plot(k, probability)
10 plt.show()
```

```
1 [1 3 4 3 3 2 2 2 2 2 1 2 4 5 3 1 2 2 2 0]
```



3.二项分布

随机变量 X 符合二项分布记为: $X \sim B(n, p)$,其中 p 是单次实验的成功率, n 是实验次数

$$E(X) = np, \sigma^2(X) = np(1 - p), P(X = k) = C_n^k p^k (1 - p)^{n-k}$$

```

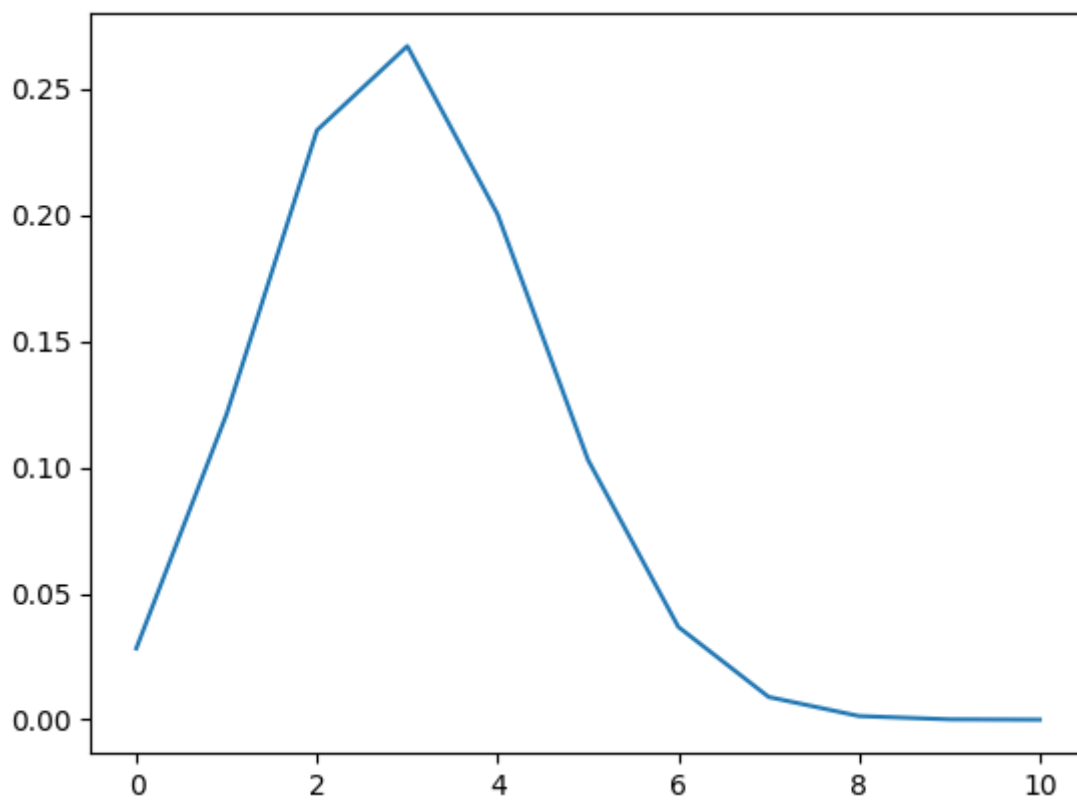
1 # 二项分布随机数生成:scipy.stats.binom.rvs(n, p, size)
2 n = 10
3 p = 0.3
4 X = sp.stats.binom.rvs(n, p, size = 10)
5 print(X)
6
7 # 二项分布概率密度函数:scipy.stats.binom.pmf(X, n, p)
8 X = np.arange(0,11)
9 Probability = sp.stats.binom.pmf(X, n, p)
10 plt.plot(X,Probability)
11 plt.show()

```

```

1 [2 4 4 3 4 6 0 2 5 3]

```



4.指数分布

随机变量 X 符合参数为 λ 或者 β 的指数分布记为: $X \sim Exp(\lambda)$ 或者 $X \sim Exp(\beta)$

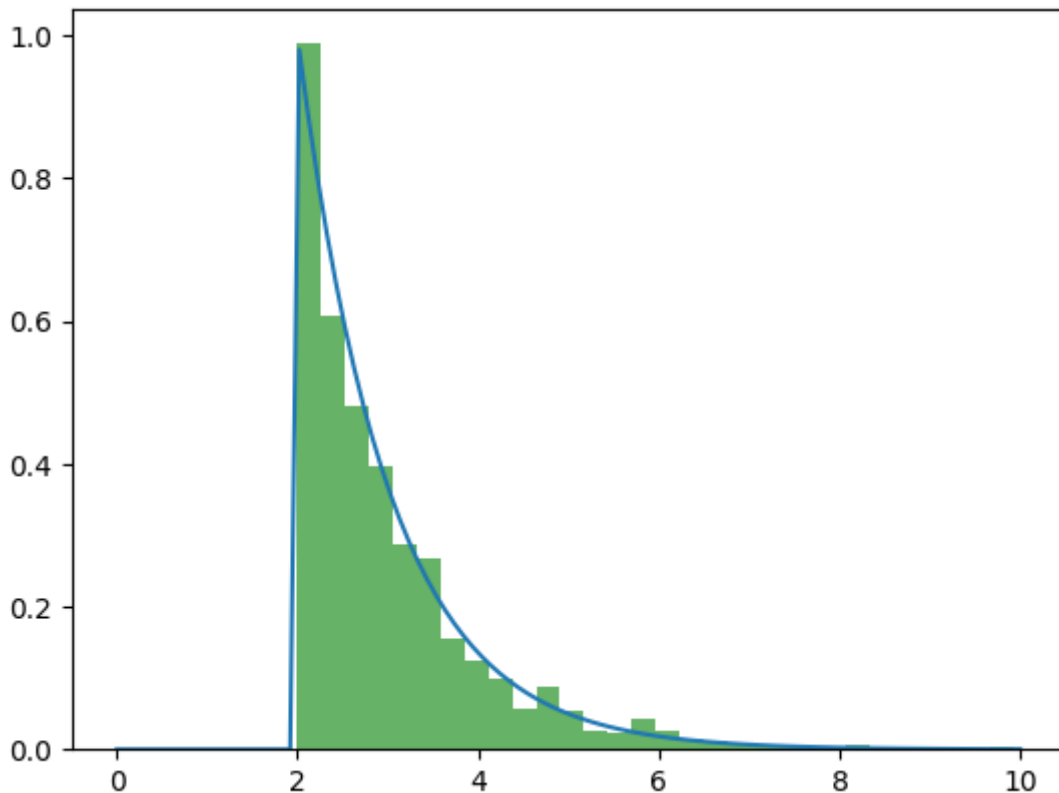
$$E(X) = \frac{1}{\lambda} = \beta, \sigma^2(X) = \frac{1}{\lambda^2} = \beta^2$$

$$f(x, \lambda) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0 \\ 0 & x < 0 \end{cases}$$

```

1 # 指数分布随机数:scipy.stats.expon.rvs(bata, size)
2 lambda_exp = 0.5
3 beta = 1/lambda_exp
4 samples = sp.stats.expon.rvs(beta, size = 1000)
5 plt.hist(samples, bins = 30, density=True, alpha=0.6, color='g',
6 label="samples")
7 # 指数分布概率函数:scipy.stats.expon.pdf(X,beta)
8 X = np.linspace(0,10,100)
9 probability = sp.stats.expon.pdf(X, beta)
10 plt.plot(X,probability)
11 plt.show()

```



5.超几何分布

随机变量X服从参数为K, N的超几何分布记为: $X \sim H(n, K, N)$

$$E(X) = \frac{nK}{N}, \sigma^2(X) = \frac{nK(N-K)(N-n)}{N^2(N-1)}$$

其概率密度函数为: $f(k; n, K, N) = \frac{C_K^k C_{N-K}^{n-k}}{C_N^n}$

```

1 # 参数定义
2 N = 100 # 100件
3 K = 30 # 次品
4 n = 20
5
6 # 超几何分布随机数:scipy.stats.hypergeom.rvs(N, K, n, size)
7 X = sp.stats.hypergeom.rvs(N, K, n, size = 10)
8 print(X)
9

```

```

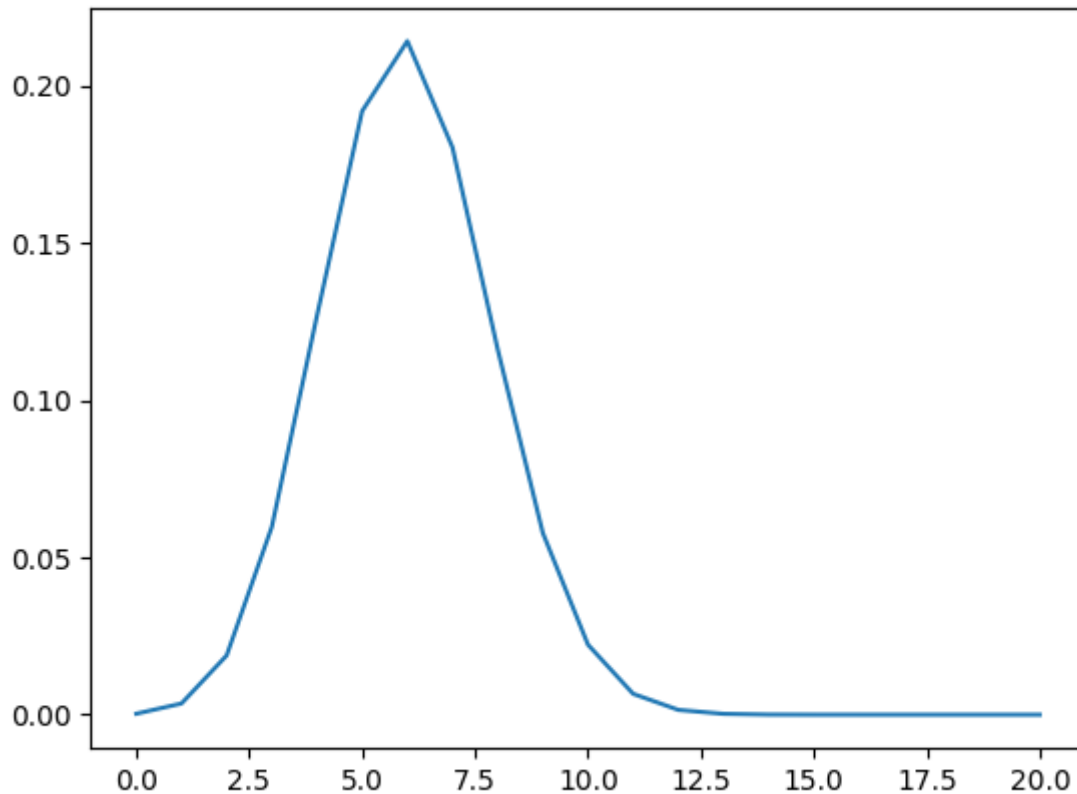
10 # 超几何分布概率函数:scipy.stats.hypergeom.pmf(k, N, K, n)
11 k = np.arange(0,n+1)
12 probability = sp.stats.hypergeom.pmf(k, N, K, n)
13 plt.plot(k, probability)
14 plt.show()

```

```

1 [5 4 8 7 6 5 4 7 8 5]

```



```

1
2
3

```

贝叶斯分类器

原理：贝叶斯分类器基于贝叶斯定理，旨在利用已知数据来模拟出新数据属于不同类别的概率，然后选择具有最高概率的类别作为分类结果，做出相应的决策

```

1 class NaiveBayesClassifier:
2     def __init__(self) :
3         self.classes = set() # 存储所有C_i类别
4         self.classes_prob = dict() # 每个C_i出现的可能性
5         self.condition_prob = dict() # 每个C_i在A下的条件概率
6         self.P0 = 1 # 先验概率
7
8     # 训练函数
9     def fit(self, spam_data_matrix, ham_data_matrix):

```

```

10     from math import log
11     spam_dim = len(spam_data_matrix)
12     ham_dim = len(ham_data_matrix)
13     dim = spam_dim + ham_dim
14     self.P0 = log(spam_dim / dim)
15
16     for line in spam_data_matrix:
17         self.classes.update(set(line))
18     for line in ham_data_matrix:
19         self.classes.update(set(line))
20
21     # C_i总体概率 和 条件概率
22     for word in self.classes:
23         spamif = 0
24         hamif = 0
25         for line in spam_data_matrix:
26             if word in line:
27                 spamif += 1
28         for row in ham_data_matrix:
29             if word in row:
30                 hamif += 1
31         self.classes_prob[word] = (spamif + hamif) / dim
32         self.condition_prob[word] = (spamif) / spam_dim
33     class_prob_sum = sum(self.classes_prob.values())
34     condition_prob_sum = sum(self.condition_prob.values())
35     self.classes_prob = {key:value/class_prob_sum for key,value in
self.classes_prob.items()}
36     self.condition_prob = {key:value/condition_prob_sum for key,value
in self.condition_prob.items()}
37
38     spam_p = (spam_dim + ham_dim) / dim
39     return spam_dim, self.classes_prob, self.condition_prob
40
41     # 判断每一个行是不是一个符合A
42     def IsSpam(self,data_list):
43         from math import log
44         condition_prob = 0
45         # 在计算条件概率时, 添加平滑值 (alpha) 到分子和分母中
46         alpha = 1 # laplace平滑参数
47         dim = len(self.classes)
48         for word in data_list:
49             if word in self.classes:
50                 condition_prob += log((self.condition_prob[word] + alpha) /
(self.classes_prob[word] + alpha * dim))
51
52         P = self.P0 + condition_prob
53
54         # 将1/2的对数概率作为阈值来判断
55         threshold = log(1/2)
56         if P >= threshold:
57             print(P)
58             return True
59         else:
60             return False
61
62     # 做总体预测

```

```

63     def predict(self, data_matrix):
64         total_dim = len(data_matrix)
65         issпам = []
66         for item in data_matrix:
67             if self.IsSpam(item) == True:
68                 issпам.append(True)
69             else:
70                 issпам.append(False)
71         return issпам
72
73
74 class BayesFile:
75     def __init__(self, pwd) :
76         self.pwd = pwd
77         self.excluded_words = set()
78         self.load_excluded_words()
79
80     def load_excluded_words(self):
81         with open(self.pwd+"/excluded_words.txt", "r", encoding="utf-8") as
file:
82             self.excluded_words = set(word.strip() for word in
file.readlines())
83
84     def process_file(self, file_path):
85         # 读取文件并将文章分割成单词
86         words = []
87         with open(file_path, 'r', encoding='utf-8') as file:
88             text = file.read()
89             words = text.split() # 简单地使用空格分割
90         words = [word for word in words if word.lower() not in
self.excluded_words]
91         return words
92
93     def build_matrix(self):
94         from os import path
95         SpamMatrix = []
96         HamMatrix = []
97         for i in range(1, 26):
98             file_name = f'{i}.txt'
99             file_path = path.join(self.pwd+"/email/spam", file_name)
100             if path.isfile(file_path):
101                 words = self.process_file(file_path)
102                 SpamMatrix.append(words)
103             file_path = path.join(self.pwd+"/email/ham", file_name)
104             if path.isfile(file_path):
105                 words = self.process_file(file_path)
106                 HamMatrix.append(words)
107
108         return SpamMatrix, HamMatrix
109
110     def detect():
111         NBC = NaiveBayesClassifier()
112         FileMatrix =
BayesFile(pwd="/home/jxluo/ubuntu/WorkPlace/grade3/Computer-
Simulations/Homeworks/hw02")
113

```

```

114     FileMatrix.load_excluded_words()
115     ham, spam = FileMatrix.build_matrix()
116
117     NBC.fit(ham_data_matrix=ham, spam_data_matrix=spam)
118     # 随机选择一些样本
119     def choice():
120         indexLst = []
121         detect_file = []
122         from random import choice
123         for index in range(20):
124             k = choice([0,1])
125             index = choice(range(1,25))
126             if k == 1:
127                 indexLst.append(True)
128                 detect_file.append(spam[index])
129             else:
130                 indexLst.append(False)
131                 detect_file.append(ham[index])
132         return detect_file, indexLst
133
134     Matrix, lst = choice()
135     result = NBC.predict(Matrix)
136     RightNum = sum(1 if result[index] == lst[index] else 0 for index in
137 range(len(lst)))
138     ratio = RightNum/len(lst)
139     return ratio
140
141 print("检测成功率:{:.2f}%".format(detect()*100))
142

```

1 | 检测成功率:65.00%

上述代码创建了一个朴素bayes选择器类NaiveBayesClassifier，用类BayesFiles做文件处理，传递给选择器。重现了实例demo_email_spam.pdf

