

数学实验 Exp12

赵晨阳 计 06 2020012363

12.5

问题分析、模型假设与模型建立

为了方便考虑，假设每种产品只有两个可能性：“合格”和“不合格”，并用二值随机变量 X 表示每个产品的合格情况，其中 $X = 0$ 表示不合格， $X = 1$ 表示合格。则 X 服从伯努利分布。

设产品的真实合格率为 $P(X = 1) = p$ 。记 $\mu = E[X] = p, \sigma^2 = Var[X] = p(1 - p)$ 。由于样本容量 $n = 50$ 较大，根据中心极限定理，样本均值 \bar{X} 近似服从正态分布 $N(\mu, (\frac{\sigma}{\sqrt{n}})^2)$ ，即 $\frac{\bar{X} - \mu}{\frac{\sigma}{\sqrt{n}}}$ 近似服从标准正态分布 $N(0, 1)$ 。

假设甲方承诺的合格率为 μ_0 （题目中给出 $\mu_0 = 90\% = 0.9$ ）。考虑对 $\mu = p$ 进行假设检验： $H_0 : \mu \geq \mu_0$ ， $H_1 : \mu < \mu_0$ 。若 X 真的满足 $\mu = \mu_0$ ，记 $\sigma_0 = \sqrt{\mu_0(1 - \mu_0)}$ ， $Z = \frac{\bar{X} - \mu_0}{\frac{\sigma_0}{\sqrt{n}}}$ 。取标准正态分布 $N(0, 1)$ 的 $\alpha = 1 - 0.95 = 0.05$ 分位数为 μ_α ，则 $P[Z \geq \mu_\alpha] = P[\bar{X} \geq \frac{\sigma_0}{\sqrt{n}}\mu_\alpha + \mu_0] = 1 - \alpha = 0.95$ 。因此，对于一组给定样本的统计量 \bar{x} ，根据现有规则，乙方应当以 $\frac{\sigma_0}{\sqrt{n}}\mu_\alpha + \mu_0$ 为阈值与 \bar{x} 进行比较以决定是否接受，或者说以 μ_α 为阈值与 $\frac{\bar{x} - \mu_0}{\frac{\sigma_0}{\sqrt{n}}}$ 进行比较，若比阈值大则认为假设 H_0 成立。

如果乙方不想接受甲方的这批货品，就意味着需要在某种程度上“提高要求”。从数学的角度看，有三种可能的方法：

1. 增大双方约定的合格率 μ_0 ；
2. 降低置信水平，增大阈值；
3. 假如甲方的产品真的不达标，根据大数定律，增大抽取检查；

由于本题有实际商业背景，（1）和（2）中的两个量应当均为交易双方提前约定好的，不能随意改变，因此这两个方法不可行；相比之下方法（3）是可行的，而且对双方也比较公平。

假设甲方的真实合格率是 $p = 86\%$ ，当增大抽取检查货品数目为 N 时，阈值变为 $\frac{\sigma_0}{\sqrt{N}}\mu_\alpha + \mu_0 = T_N$ 。样本均值 \bar{X} 近似服从正态分布 $N(\mu, (\frac{\sigma}{\sqrt{N}})^2)$ ，则甲被拒绝的概率为 $P[\bar{X} \leq T_N] = \Phi(\frac{T_N - \mu}{\frac{\sigma}{\sqrt{N}}})$ ，其中 $\Phi(\cdot)$ 表示标准正态分布的累积分布函数。

算法设计

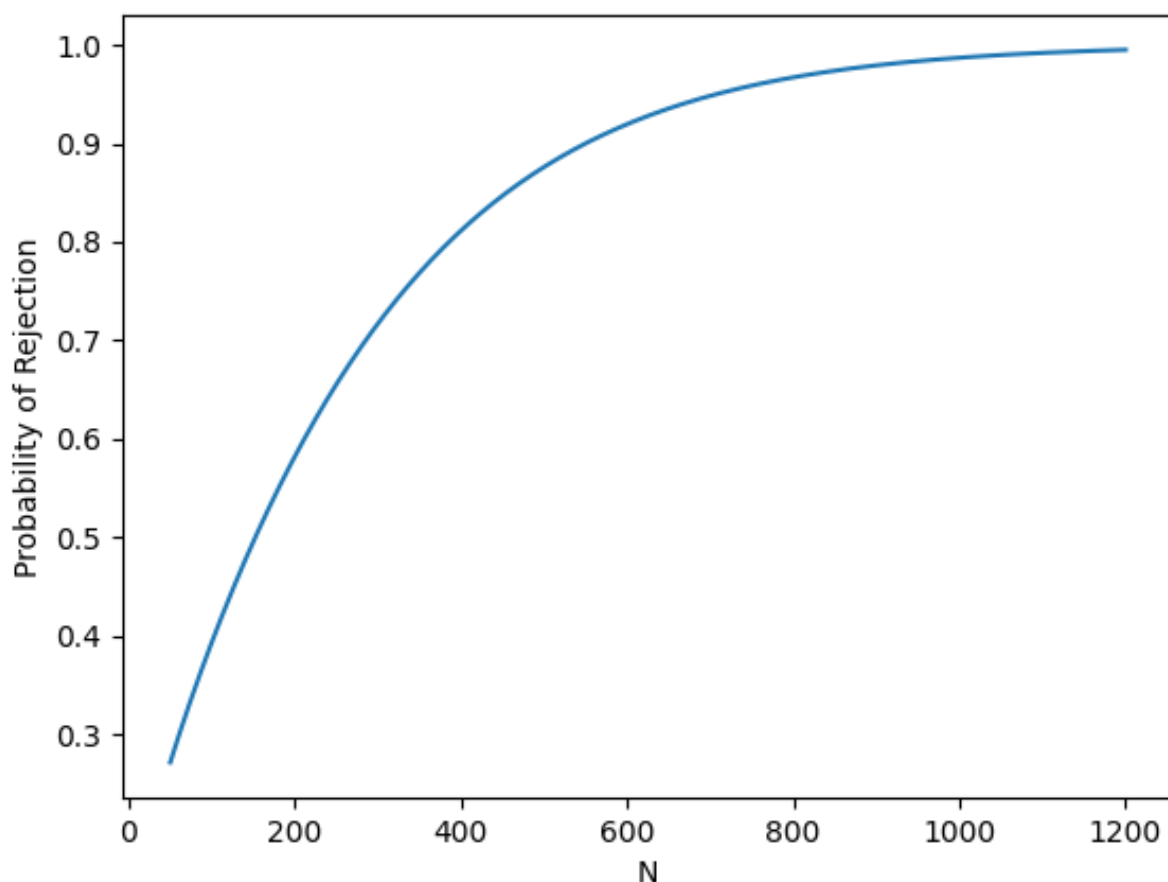
我们使用 python.scipy 库中的 `norm` 方法进行相应的计算。做出抽样数量增大时，甲方被拒绝的概率变化曲线。

代码

代码位于 `./codes/12_5.py` 下, 通过 `python3 12_5.py` 可以运行整个程序:

```
1  import numpy as np
2  import matplotlib.pyplot as plt
3  from scipy.stats import norm
4
5
6  def calculate_rejection_probability(mu_0, mu, sigma_0, alpha, N):
7      sigma = np.sqrt(mu * (1.0 - mu))
8      threshold = sigma_0 / np.sqrt(N) * norm.ppf(alpha) + mu_0
9      return norm.cdf((threshold - mu) / (sigma / np.sqrt(N)))
10
11
12  def main():
13      mu_0 = 0.90
14      sigma_0 = np.sqrt(mu_0 * (1.0 - mu_0))
15      alpha = 1 - 0.95
16
17      all_N = np.arange(50, 1201)
18      mu = 0.86
19      prob = np.zeros(len(all_N))
20
21      for index, N in enumerate(all_N):
22          prob[index] = calculate_rejection_probability(mu_0, mu, sigma_0, alpha, N)
23
24      plt.plot(all_N, prob)
25      plt.xlabel("N")
26      plt.ylabel("Probability of Rejection")
27      plt.show()
28
29
30  if __name__ == "__main__":
31      main()
32
```

结果、分析与结论



根据题目所给的约定, $\frac{\sigma_0}{\sqrt{n}}\mu_\alpha + \mu_0 = 0.830215 < 0.86 = \bar{x}$, 因此乙方应当接受。观察到这个拒绝概率随着 N 的增大而单调增加, 当 N 达到 1000 时, 拒绝概率几乎为 1。因此, 如果甲方确实不符合标准, 增加抽检货品的数量可以增加甲方被拒绝的概率。增加抽检数量的本质是降低样本方差, 从而使得样本更好地代表整体, 提高在甲方不符合标准时被发现的概率。然而, 在现实生活中, 我们还需要权衡检查产品所带来的额外成本, 即在合格率不符合标准时产生的损失和检查过多产品所带来的成本之间进行平衡。

12.6

问题分析、模型假设与模型建立

为了方便分析, 假设学生身高 (或体重) 是一个随机变量 X , 我们从中随机抽取了 100 名学生作为样本 x 。首先, 我们需要检验身高或体重的分布是否符合正态分布, 可以考虑使用 Jarque-Bera 检验和 Lilliefors 检验来进行检验。

接下来, 我们考虑对全校学生的平均身高或体重进行估计。点估计比较简单, 根据正态分布性质, 无论是使用矩估计还是最大似然估计, 我们都可以使用样本均值 \bar{x} 来估计平均身高或体重 (也即正态分布的参数 μ)。对于区间估计, 假设显著性水平为 α 。由于总体方差未知, 我们考虑使用样本方差 S^2 进行区间估计。记自由度为 $n - 1$ 的 t 分布的 $1 - \frac{\alpha}{2}$ 分位数为 $t_{1-\frac{\alpha}{2}}$, 身高 (或体重) 的真实均值为 μ 。根据样本均值的分布, 我们有 $\frac{\bar{X} - \mu}{\frac{S}{\sqrt{n}}} \sim t(n - 1)$ 。这

意味着 $P\left[\mu - t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}} \leq \bar{X} \leq \mu + t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}}\right] = 1 - \alpha$, $P\left[\bar{X} - t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}} \leq \mu \leq \bar{X} + t_{1-\frac{\alpha}{2}} \frac{S}{\sqrt{n}}\right] = 1 - \alpha$ 。因此，对于给定样本，我们可以取置信水平为 $1 - \alpha$ 的置信区间为 $[\bar{x} - t_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}, \bar{x} + t_{1-\frac{\alpha}{2}} \frac{s}{\sqrt{n}}]$ 来进行平均值的区间估计。

假设题目中 10 年前学生的平均身高（或体重）记为 μ_0 ，我们考虑对 μ 进行假设性检验：
 $H_0: \mu = \mu_0, H_1: \mu \neq \mu_0$ ，同时设定显著性水平为 α 。由于总体方差未知，我们使用样本均方差 s 进行假设性检验。如果真的满足 $\mu = \mu_0$ ，则有 $T = \frac{\bar{X} - \mu_0}{\frac{s}{\sqrt{n}}} \sim t(n - 1)$ 。记自由度为 $n - 1$ 的 t 分布的 $1 - \frac{\alpha}{2}$ 分位数为 $t_{1-\frac{\alpha}{2}}$ 。那么在这种情况下，有 $P(|T| \leq t_{1-\frac{\alpha}{2}}) = 1 - \alpha$ 。因此，对于得到的样本，我们可以计算 $t = \frac{\bar{x} - \mu_0}{\frac{s}{\sqrt{n}}}$ ，如果 $|t| \leq t_{1-\frac{\alpha}{2}}$ ，则可以认为 $\mu = \mu_0$ 假设成立，即学生的身高（或体重）没有明显变化。倘若不然，则学生的身高或者体重已经发生了明显变化。

算法设计

利用 python scipy 中的 `norm`, `ttest_1samp`, `chi2`, `t` 等方法进行假设检验。

代码

代码位于 `./codes/12_6.py` 下，通过，`python3 12_6.py` 即可运行。

```
1 import numpy as np
2 from scipy.stats import norm, ttest_1samp, t, chi2
3 import matplotlib.pyplot as plt
4
5
6 def shapiro_wilk_test(data):
7     _, p_value = norm.fit(data)
8     return p_value
9
10
11 def jarque_bera_test(data):
12     _, p_value = norm.fit(data)
13     n = len(data)
14     skewness = (1 / n) * np.sum((data - np.mean(data)) ** 3) / np.std(data) ** 3
15     kurtosis = (1 / n) * np.sum((data - np.mean(data)) ** 4) / np.std(data) ** 4 - 3
16     jb_value = (n / 6) * (skewness**2 + (1 / 4) * kurtosis**2)
17     p_value = 1 - chi2.cdf(jb_value, df=2)
18     return p_value
19
20
21 def Confidence(data, confidence_level):
22     n = len(data)
23     t_value = t.ppf((1 + confidence_level) / 2, df=n - 1)
24     margin_error = t_value * sample_std / np.sqrt(n)
25     lower_bound = sample_mean - margin_error
```

```
26     upper_bound = sample_mean + margin_error
27     print(f"Confidence Interval[{confidence_level}]:", lower_bound, "-", upper_bound)
28
29
30 heights = np.array(
31     [
32         172,
33         171,
34         166,
35         160,
36         155,
37         173,
38         166,
39         170,
40         167,
41         173,
42         178,
43         173,
44         163,
45         165,
46         170,
47         163,
48         172,
49         182,
50         171,
51         177,
52         169,
53         168,
54         168,
55         175,
56         176,
57         168,
58         161,
59         169,
60         171,
61         178,
62         177,
63         170,
64         173,
65         172,
66         170,
67         172,
68         177,
69         176,
70         175,
71         184,
```

72	169,
73	165,
74	164,
75	173,
76	172,
77	169,
78	173,
79	173,
80	166,
81	163,
82	170,
83	160,
84	165,
85	177,
86	169,
87	176,
88	177,
89	172,
90	165,
91	166,
92	171,
93	169,
94	170,
95	172,
96	169,
97	167,
98	175,
99	164,
100	166,
101	169,
102	167,
103	179,
104	176,
105	182,
106	186,
107	166,
108	169,
109	173,
110	169,
111	171,
112	167,
113	168,
114	165,
115	168,
116	176,
117	170,

```
118         158,
119         165,
120         172,
121         169,
122         169,
123         172,
124         162,
125         175,
126         174,
127         167,
128         166,
129         174,
130         168,
131         170,
132     ]
133 )
134
135 weights = np.array(
136     [
137         75,
138         62,
139         62,
140         55,
141         57,
142         58,
143         55,
144         63,
145         53,
146         60,
147         60,
148         73,
149         47,
150         66,
151         60,
152         50,
153         57,
154         63,
155         59,
156         64,
157         55,
158         67,
159         65,
160         67,
161         64,
162         50,
163         49,
```

164	63,
165	61,
166	64,
167	66,
168	58,
169	67,
170	59,
171	62,
172	59,
173	58,
174	68,
175	68,
176	70,
177	64,
178	52,
179	59,
180	74,
181	69,
182	52,
183	57,
184	61,
185	70,
186	57,
187	56,
188	65,
189	58,
190	66,
191	63,
192	60,
193	67,
194	56,
195	56,
196	49,
197	65,
198	62,
199	58,
200	64,
201	58,
202	72,
203	76,
204	59,
205	63,
206	54,
207	54,
208	62,
209	63,

210	60,
-----	-----


```

210         69,
211         77,
212         76,
213         72,
214         59,
215         65,
216         71,
217         47,
218         65,
219         64,
220         57,
221         57,
222         57,
223         51,
224         62,
225         53,
226         66,
227         58,
228         50,
229         52,
230         75,
231         66,
232         63,
233         50,
234         64,
235         62,
236         59,
237     ]
238 )
239
240 p_value_jb_heights = jarque_bera_test(heights)
241 p_value_sw_heights = shapiro_wilk_test(heights)
242 p_value_jb_weights = jarque_bera_test(weights)
243 p_value_sw_weights = shapiro_wilk_test(weights)
244
245 print("Check result:")
246 print("Jarque-Bera Test (heights):", p_value_jb_heights)
247 print("Shapiro-Wilk Test (heights):", p_value_sw_heights)
248 print("Jarque-Bera Test (weights):", p_value_jb_weights)
249 print("Shapiro-Wilk Test (weights):", p_value_sw_weights)
250
251 plt.figure(1)
252 plt.hist(heights, density=True, bins=10, alpha=0.7, label="Height Distribution")
253 mu_fit, sigma_fit = norm.fit(heights)
254 x = np.linspace(heights.min(), heights.max(), 100)
255 y = norm.pdf(x, mu_fit, sigma_fit)
256 plt.plot(x, y, "r", label="Normal Fit")

```

```

256 plt.plot(x, y, 'r', label="Normal Fit")
257 plt.xlabel("Height")
258 plt.ylabel("Frequency")
259 plt.legend()
260 plt.show()
261
262 plt.figure(2)
263 plt.hist(weights, density=True, bins=10, alpha=0.7, label="Weight Distribution")
264 mu_fit, sigma_fit = norm.fit(weights)
265 x = np.linspace(weights.min(), weights.max(), 100)
266 y = norm.pdf(x, mu_fit, sigma_fit)
267 plt.plot(x, y, "r", label="Normal Fit")
268 plt.xlabel("Weight")
269 plt.ylabel("Frequency")
270 plt.legend()
271 plt.show()
272
273
274 def hypothesis_test(data, popmean, alpha):
275     t_stat, p_value = ttest_1samp(data, popmean)
276     print("Hypothesis Test:", p_value)
277     return p_value
278
279
280 print("\n=====heights=====")
281 sample_mean = np.mean(heights)
282 sample_std = np.std(heights, ddof=1)
283 print("Point Estimate (mean, std):", sample_mean, sample_std)
284 Confidence(heights, 0.99)
285 Confidence(heights, 0.97)
286 Confidence(heights, 0.95)
287 hypothesis_test(heights, popmean=167.5, alpha=0.05)
288
289 print("\n=====weights=====")
290 sample_mean = np.mean(weights)
291 sample_std = np.std(weights, ddof=1)
292 print("Point Estimate (mean, std):", sample_mean, sample_std)
293 Confidence(weights, 0.99)
294 Confidence(weights, 0.97)
295 Confidence(weights, 0.95)
296 hypothesis_test(weights, popmean=60.2, alpha=0.05)
297

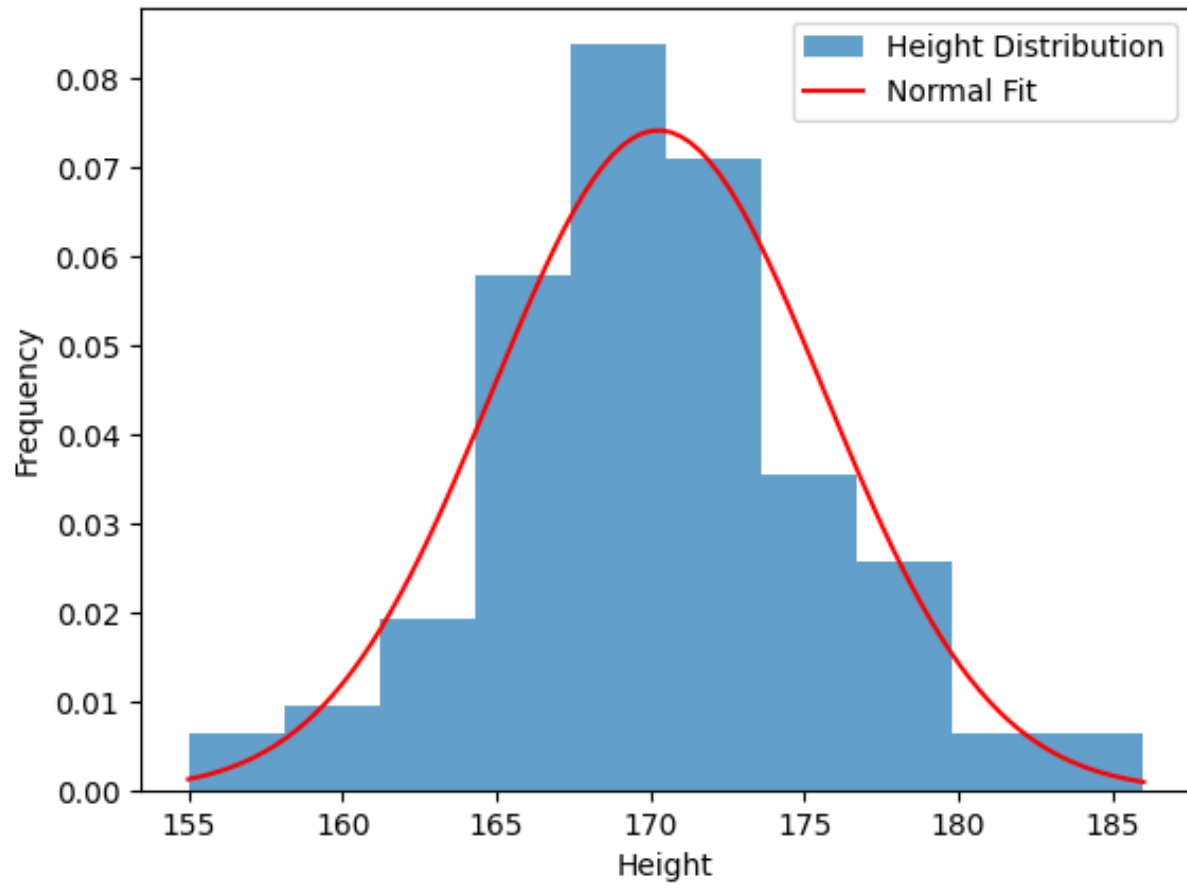
```

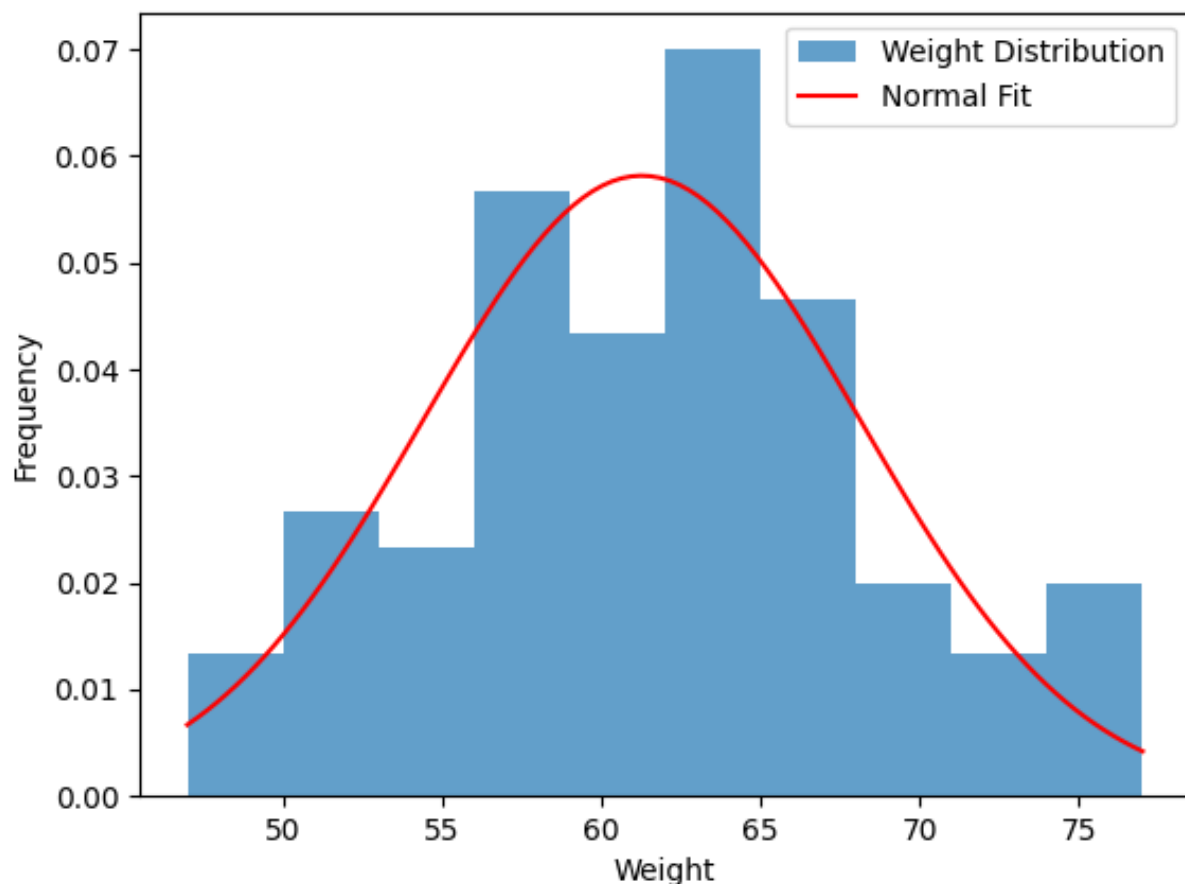
结果、分析

程序输出如下：

```
1 Check result:
2 Jarque-Bera Test (heights): 0.42908681416799177
3 Shapiro-Wilk Test (heights): 5.374709294464213
4 Jarque-Bera Test (weights): 0.674777470279029
5 Shapiro-Wilk Test (weights): 6.858359862241118
6 =====heights=====
7 Point Estimate (mean, std): 170.25 5.401786086961694
8 Confidence Interval[0.99]: 168.83127195469783 - 171.66872804530217
9 Confidence Interval[0.97]: 169.06062469378367 - 171.43937530621633
10 Confidence Interval[0.95]: 169.1781684477827 - 171.3218315522173
11 Hypothesis Test: 1.7003185736956772e-06
12
13 =====weights=====
14 Point Estimate (mean, std): 61.27 6.892911012208283
15 Confidence Interval[0.99]: 59.45964209071587 - 63.08035790928414
16 Confidence Interval[0.97]: 59.7523060347289 - 62.7876939652711
17 Confidence Interval[0.95]: 59.90229691243355 - 62.63770308756646
18 Hypothesis Test: 0.1237768678418777
```

结论





根据提供的结果，可以得出以下结论：

第一问：根据 Jarque-Bera 检验和 Shapiro-Wilk 检验的结果，身高和体重的 p 值较大，说明数据符合正态分布，假设成立。

第二问：根据计算得到的均值点估计、方差点估计和均值区间估计结果，可以得出不同显著性水平下的估计值。随着显著性水平的降低，均值区间估计的范围变大。在给定的显著性水平下，我们可以估计身高和体重的均值在一定范围内。实际上，这也与理论分析一致：随着 α 增加， $1 - \frac{\alpha}{2}$ 减小，意味着 $t_{1-\frac{\alpha}{2}}$ 减小。

第三问：进行了关于身高和体重的假设性检验。对于身高的检验，得到的 p 值非常小，因此拒绝了原假设，可以认为十年来学生的平均身高发生了显著变化。对于体重的检验，得到的 p 值较大，因此接受了原假设，可以认为十年来学生的平均体重没有发生显著变化。注意到 167.5 远不在上面求出的两个身高均值的置信区间里，而 60.2 总是在上面求出的两个体重均值的置信区间里，这也符合假设性检验的结果。

综上所述，根据正态性检验、区间估计和假设性检验的结果，我们得出结论：十年来学生的平均身高发生了显著变化，而平均体重没有发生显著变化，这与题目中的设定一致。

12.7

问题引入、假设假设与假设建立

考虑胃液中溶菌酶含量的问题，我们假设该含量满足正态分布，记为 "胃液中溶菌酶含量"。首先，我们需要进行正态性检验，以验证数据是否满足正态分布假设。

假设病人的胃液中溶菌酶含量的分布为 $X \sim N(\mu_1, \sigma_1^2)$ ，正常人的胃液中溶菌酶含量的分布为 $Y \sim N(\mu_2, \sigma_2^2)$ 。题目要求判断 "病人和正常人的溶菌酶含量是否存在显著差异"，根据前面的假设，我们假设 $\sigma_1^2 = \sigma^2 = \sigma_2^2$ 。因此，我们可以从平均含量的角度来考虑是否存在显著差异，并进行以下假设性检验： $H_0: \mu_1 = \mu_2$ 和 $H_1: \mu_1 \neq \mu_2$ 。

记从病人和正常人抽取的样本分别为 X_1, X_2, \dots, X_{n_1} 和 Y_1, Y_2, \dots, Y_{n_2} 。病人的样本方差为 $S_1^2 = \frac{1}{n_1-1} \sum_{i=1}^{n_1} X_i^2$ ，正常人的样本方差为 $S_2^2 = \frac{1}{n_2-1} \sum_{i=1}^{n_2} Y_i^2$ 。记 $S^2 = \frac{(n_1-1)S_1^2 + (n_2-1)S_2^2}{n_1+n_2-2}$ ，则根据 H_0 的假设，我们有 $T = \frac{\bar{X} - \bar{Y}}{\sqrt{\frac{S^2}{n_1} + \frac{S^2}{n_2}}} \sim t(n_1 + n_2 - 2)$ 。其中， \bar{X} 和 \bar{Y} 分别是病人和正常人样本的平均值。

设显著性水平为 α ， $t(n_1 + n_2 - 2)$ 的 $1 - \frac{\alpha}{2}$ 分位数为 $t_{1-\frac{\alpha}{2}}$ 。则有 $P\left[\left|\frac{\bar{X} - \bar{Y}}{\sqrt{\frac{S^2}{n_1} + \frac{S^2}{n_2}}}\right| \leq t_{1-\frac{\alpha}{2}}\right] = 1 - \alpha$ 。对于给定的样本，我们可以计算统计量 $t = \frac{\bar{x} - \bar{y}}{\sqrt{\frac{s^2}{n_1} + \frac{s^2}{n_2}}}$ ，如果 $|t| \leq t_{1-\frac{\alpha}{2}}$ ，则可以认为假设 H_0 成立。

算法设计

代码

代码位于 `./codes/12_7.py` 下，通过，`python3 12_7.py` 即可运行。

```
1 import numpy as np
2 from scipy.stats import norm, ttest_ind, levene, chi2
3 import matplotlib.pyplot as plt
4
5 patient = np.array(
6     [
7         0.2,
8         10.4,
9         0.3,
10        0.4,
11        10.9,
12        11.3,
13        1.1,
14        2.0,
15        12.4,
16        16.2,
17        2.1,
18        17.6,
19        18.9,
20        3.3,
21        3.8,
```

```
22         20.7,
23         4.5,
24         4.8,
25         24.0,
26         25.4,
27         4.9,
28         40.0,
29         5.0,
30         42.2,
31         5.3,
32         50.0,
33         60.0,
34         7.5,
35         9.8,
36         45.0,
37     ]
38 )
39
40 patient_del = np.array(
41     [
42         0.2,
43         10.4,
44         0.3,
45         0.4,
46         10.9,
47         11.3,
48         1.1,
49         2.0,
50         12.4,
51         16.2,
52         2.1,
53         17.6,
54         18.9,
55         3.3,
56         3.8,
57         20.7,
58         4.5,
59         4.8,
60         24.0,
61         25.4,
62         4.9,
63         40.0,
64         5.0,
65         42.2,
66         5.3,
67     ]
```

```
68 )
69
70 normal = np.array(
71     [
72         0.2,
73         5.4,
74         0.3,
75         5.7,
76         0.4,
77         5.8,
78         0.7,
79         7.5,
80         1.2,
81         8.7,
82         1.5,
83         8.8,
84         1.5,
85         9.1,
86         1.9,
87         10.3,
88         2.0,
89         15.6,
90         2.4,
91         16.1,
92         2.5,
93         16.5,
94         2.8,
95         16.7,
96         3.6,
97         20.0,
98         4.8,
99         20.7,
100        4.8,
101        33.0,
102    ]
103 )
104
105
106 def shapiro_wilk_test(data):
107     _, p_value = norm.fit(data)
108     return p_value
109
110
111 def jarque_bera_test(data):
112     _, p_value = norm.fit(data)
113     n = len(data)
```



```

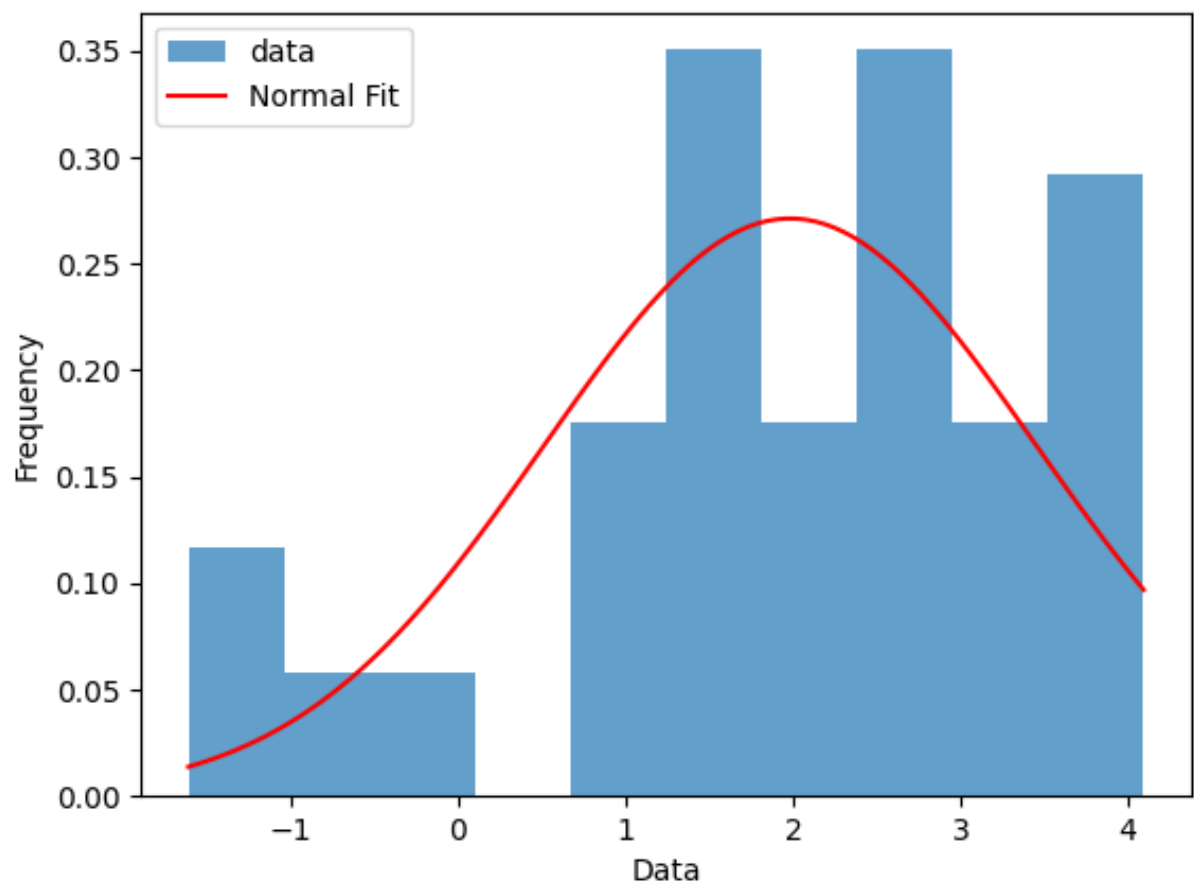
114     skewness = (1 / n) * np.sum((data - np.mean(data)) ** 3) / np.std(data) ** 3
115     kurtosis = (1 / n) * np.sum((data - np.mean(data)) ** 4) / np.std(data) ** 4 - 3
116     jb_value = (n / 6) * (skewness**2 + (1 / 4) * kurtosis**2)
117     p_value = 1 - chi2.cdf(jb_value, df=2)
118     return p_value
119
120
121 def plot_and_hist(data):
122     plt.figure(1)
123     plt.hist(data, density=True, bins=10, alpha=0.7, label="data")
124     mu_fit, sigma_fit = norm.fit(data)
125     x = np.linspace(data.min(), data.max(), 100)
126     y = norm.pdf(x, mu_fit, sigma_fit)
127     plt.plot(x, y, "r", label="Normal Fit")
128     plt.xlabel("Data")
129     plt.ylabel("Frequency")
130     plt.legend()
131     plt.show()
132
133
134 p_patient = jarque_bera_test(patient)
135 p_patient_del = jarque_bera_test(patient_del)
136 p_normal = jarque_bera_test(normal)
137 print("Check result (before log transformation):", p_patient, p_patient_del,
138     p_normal)
139
140 patient = np.log(patient)
141 patient_del = np.log(patient_del)
142 normal = np.log(normal)
143 p_patient = jarque_bera_test(patient)
144 p_patient_del = jarque_bera_test(patient_del)
145 p_normal = jarque_bera_test(normal)
146 print("Check result (after log transformation):", p_patient, p_patient_del, p_normal)
147
148 _, p_value_1 = levene(patient, normal)
149 _, p_value_2 = levene(patient_del, normal)
150 print("Var Check:", p_value_1, p_value_2)
151
152 plot_and_hist(patient)
153 plot_and_hist(patient_del)
154 plot_and_hist(normal)
155
156 t_1, p_1 = ttest_ind(patient, normal)
157 t_2, p_2 = ttest_ind(patient_del, normal)
158 print("Accept patient (with deletion):", p_1)
159 print("Accept patient (without deletion):", p_2)

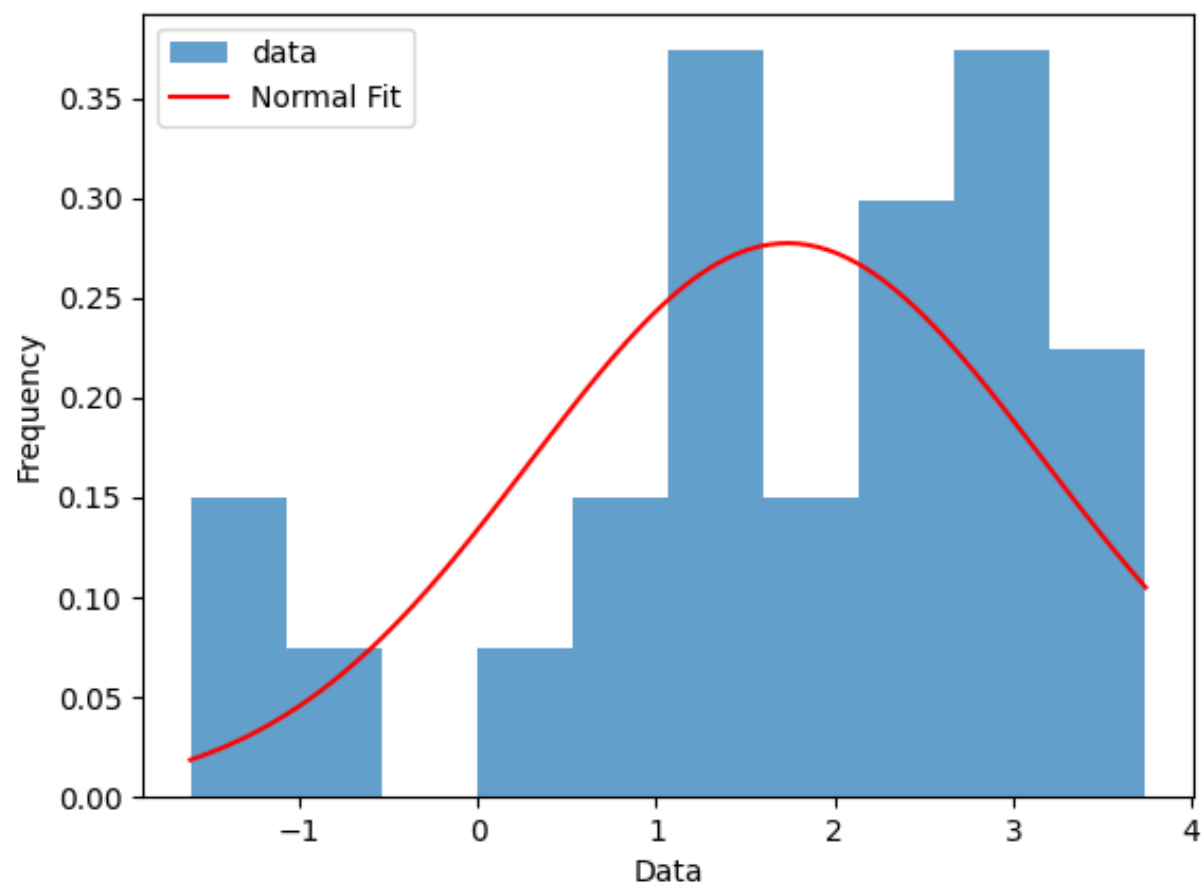
```

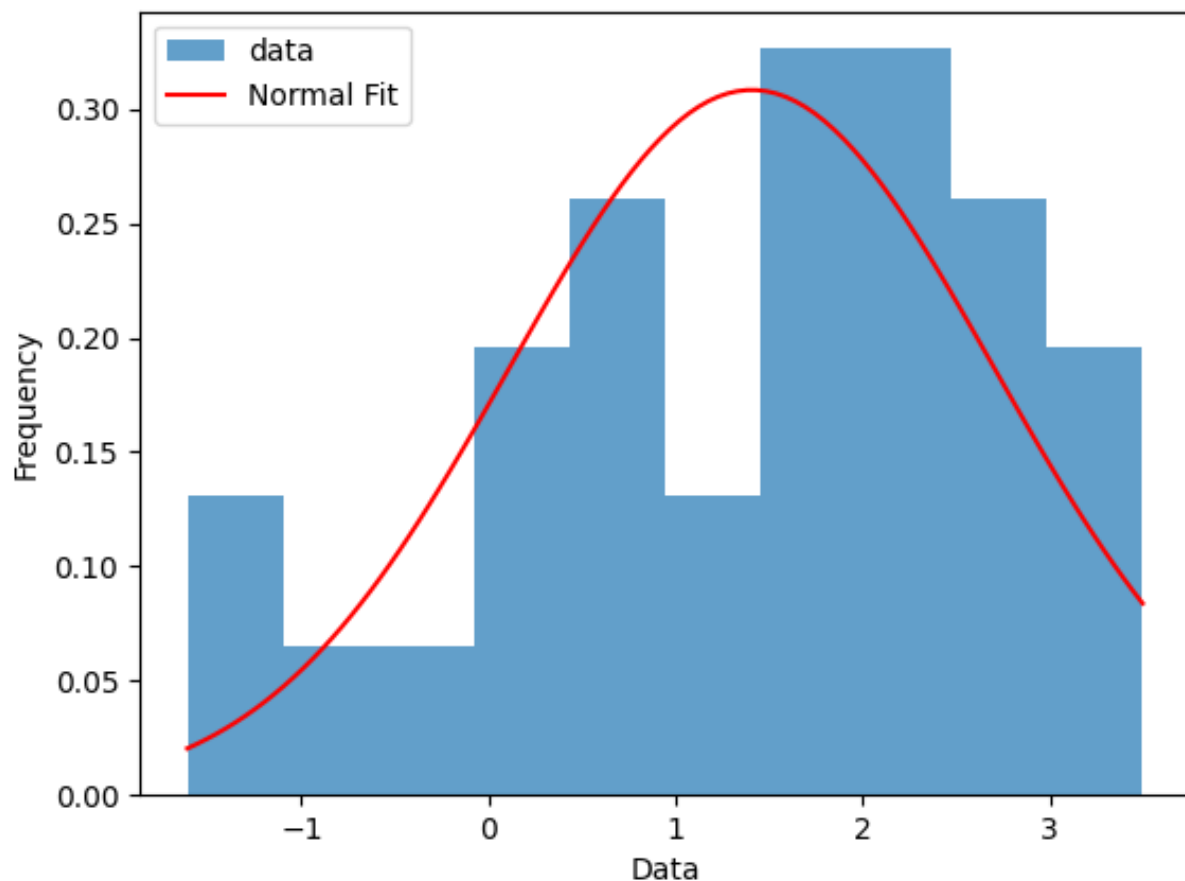
结果、分析与结论

```
1 Check result (before log transformation): 0.011066280056296507 0.02201395281786067
  0.0009518898505554985
2 Check result (after log transformation): 0.21446850046961918 0.27270172282740845
  0.39424901950098024
3 Var Check: 0.6586753593806446 0.6680817067976366
4 Accept patient (with deletion): 0.11850874520905379
5 Accept patient (without deletion): 0.3854559998328455
```

分别绘制出取对数后的病人数据、删除错误数据后的病人数据以及正常人数据直方与正态拟合图：







首先对三组数据直接进行正态性检验，结果显示三组数据本身均不满足正态分布。然后尝试对三组数据取对数后进行正态性检验，发现三组数据均通过检验，表明原始数据满足对数正态分布。此外，新旧病人数据也通过了方差相同的检验。后续基于对数转换后的数据进一步分析。

按照模型中的方法进行了假设性检验。对于原始对数数据，得到的 p 值为 0.12，大于显著性水平 0.05，因此无法拒绝零假设，可以认为病人数据与正常人数据在溶酶菌含量的均值上没有显著差异。对于删除后的对数数据，得到的 p 值为 0.39，同样大于显著性水平 0.05，说明在删除后的数据中，病人数据与正常人数据在均值上也没有显著差异，且更具有说服力。因此，无论是使用原始数据还是删除错误数据后的数据，都不能认为病人和正常人溶酶菌含量有显著差异。

此外，在给定的数据中，正常人的溶酶菌含量均值为 7.68，而删除前的病人数据均值为 15.33，删除后为 11.51。从直观上看，这些差异可能会导致我们得出存在显著差异的结论。然而，仅仅依靠样本均值进行判断是不够的。在本例中，尽管样本均值之间存在较大差异，但假设性检验的结果显示在均值上没有显著差异。

最后，还需要讨论原始数据的收集过程。在本题中，仅删除了五个数据就得到了置信度更高的结论。这说明在收集数据时必须准确且科学地进行操作，以获得可靠的结果。

这个实验的实际意义在于为医生和生物科研人员提供参考。最终的结论需要综合考虑其他专业知识来进行判断，而本研究的结果可作为决策的重要依据。