

第一次仿真实验报告

冯海桐-522031910557

2024 年 11 月 10 日

1 仿真实验题1

1.1 实验描述

通过仿真实验验证高维数据的集中分布。利用高斯分布随机函数生成维数 d ($d = 1000, 500, 100, 3$) 四组样本数据 (样本个数 $N = 10000$)，分别计算每一个样本与原点的距离，估计每一组样本数据的其距离的分布函数，并画出相应分布函数的波形曲线。

1.2 实验分析

设高斯分布随机函数为 $X \sim N(0, 1)$ ，则 d 维高斯分布随机函数为 $X \sim N(0, I_d)$ ，其中 I_d 为 d 维单位矩阵。设 $X = (X_1, X_2, \dots, X_d)$ ， $Y = |X|$ ， $Z = Y^2$ ，则 Z 服从自由度为 d 的卡方分布，即 $Z \sim \chi^2(d)$ 。

Z 的概率密度函数为

$$f_Z(z) = \frac{1}{2^{d/2}\Gamma(d/2)} z^{d/2-1} e^{-z/2} \quad (1)$$

则 Z 的分布函数为

$$F_Z(z) = \int_0^z f_Z(t) dt = 1 - \frac{\gamma(d/2, z/2)}{\Gamma(d/2)} \quad (2)$$

其中 $\gamma(a, x) = \int_0^x t^{a-1} e^{-t} dt$ 为不完全伽马函数。

则 Y 的分布函数为

$$F_Y(y) = F_Z(y^2) = 1 - \frac{\gamma(d/2, y^2/2)}{\Gamma(d/2)} \quad (3)$$

则 Y 的概率密度函数为

$$f_Y(y) = \frac{dF_Y(y)}{dy} = \frac{y^{d-1} e^{-y^2/2}}{2^{d/2-1}\Gamma(d/2)} \quad (4)$$

其概率密度最大点为 $y = \sqrt{d-1}$ 。

1.3 实验代码

利用python编写代码如下：

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from scipy.special import gammaln
4
5 # Set the dimensions and the number of samples
```

```
6 dimensions = [1000, 500, 100, 3]
7 N = 10000
8
9 plt.figure(figsize=(12, 8))
10
11 for d in dimensions:
12     # Generate random samples
13     data = np.random.randn(N, d)
14
15     # Calculate the distances
16     distances = np.linalg.norm(data, axis=1)
17
18     # Plot the histogram of the distances
19     plt.hist(distances, bins=100, density=True, alpha=0.6, label=f'd
20               ={d}')
21
22     # Plot the theoretical distribution
23     y = np.linspace(1e-10, np.max(distances), 1000)
24     log_pdf = (d-1) * np.log(y) - y**2 / 2 - (d/2 - 1) * np.log(2) -
25               gammaln(d/2)
26     pdf = np.exp(log_pdf)
27     plt.plot(y, pdf, linewidth=2)
28
29 # Show the plot
30 plt.title('Distribution of distances from the origin')
31 plt.xlabel('Distances')
32 plt.ylabel('Density')
33 plt.legend()
34 plt.show()
```

先生成 d 维高斯分布随机函数，计算每一个样本与原点的距离，然后画出相应分布函数的波形曲线。由于 $d = 1000, 500$ 时，指数计算过于庞大，故取对数计算，最后再取指数。

1.4 实验结果

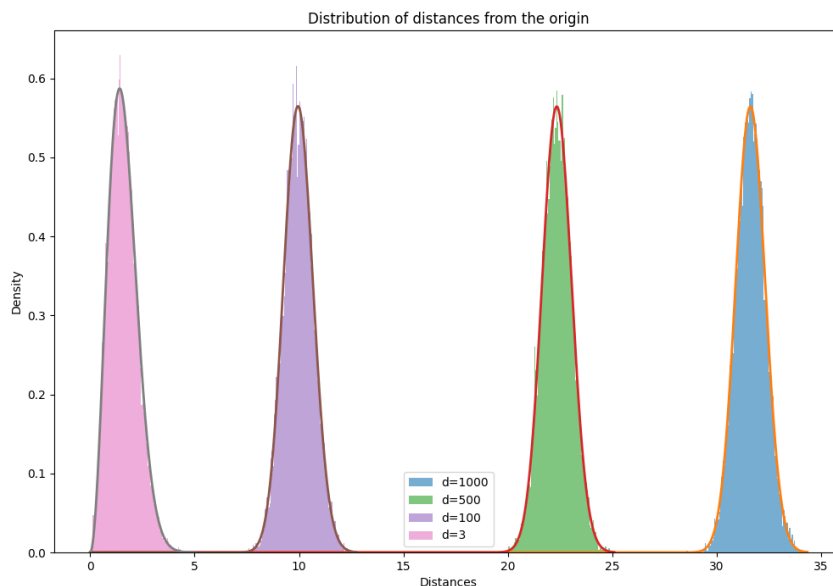


图 1: 高维数据的集中分布

由此可见，仿真实验结果与理论分析基本吻合。

2 仿真实验题2

2.1 实验描述

设 $f(x) = x_1x_2x_3 + x_3^3$ 是一个三个自变量的函数。考虑统计模型 $y = f(x) + \varepsilon, x \in R^n$ ，其中 ε 服从零均值，单位方差的高斯分布。分别对不同维数自变量 $n = 5, 10, 20, 30$ ，进行采样 $N = 1000$ 个样本数据。利用三次多项式进行回归，分别计算不同维数回归模型的方差，并讨论模型方差随着自变量维数增加时的变化规律。

2.2 实验分析

由于 $f(x) = x_1x_2x_3 + x_3^3$ 是一个三个自变量的函数，随着自变量维数的增加，回归模型的方差理应不变。但是由于采样数据的误差，回归模型随着自变量维数（不相关特征）的增加，容易产生过拟合现象，导致模型方差增大。

2.3 实验代码

利用python编写代码如下：

```
1 import numpy as np
2 from sklearn.preprocessing import PolynomialFeatures
```

```
3 from sklearn.linear_model import LinearRegression
4 from sklearn.model_selection import train_test_split
5
6 # Set the dimensions and the number of samples
7 dimensions = [5, 10, 20, 30]
8 N = 1000
9
10 # Store the variances
11 variances = []
12
13 for d in dimensions:
14     X = np.random.randn(N, d)
15     y = X[:, 0] * X[:, 1] * X[:, 2] + X[:, 2]**3 + np.random.randn(N)
16
17     # Split the data into training and testing sets
18     X_train, X_test, y_train, y_test = train_test_split(X, y,
19                                                         test_size=0.3, random_state=42)
20
21     # Fit the polynomial regression model
22     poly = PolynomialFeatures(degree=3)
23     X_train_poly = poly.fit_transform(X_train)
24     X_test_poly = poly.transform(X_test)
25
26     model = LinearRegression()
27     model.fit(X_train_poly, y_train)
28
29     # Calculate the variance
30     y_pred = model.predict(X_test_poly)
31     variance = np.var(y_pred - y_test)
32     variances.append(variance)
33
34 # Print the variances
35 for d, var in zip(dimensions, variances):
36     print(f" Model variance in dimension {d}: {var}")
```

先生成 n 维高斯分布随机函数，然后进行多项式回归，计算模型方差。

2.4 实验结果

进行三次实验，得到不同维度下的模型方差如表1所示。

维度	第一次	第二次	第三次
5	1.126	1.100	1.053
10	2.058	2.190	2.326
20	3.086	3.007	3.226
30	3.183	3.287	4.576

表 1: 不同维度下的模型方差

由此可见，随着自变量维数增加，模型出现过拟合现象，模型方差增大。