

## 3.7 方差

---

方差是各个数据与其[算术平均数](#)的[离差](#)平方和的平均数。

**方差：** 
$$S^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n - 1}$$

除以n-1而不是n，是因为这样能使我们以较小的样本集更好地逼近总体的标准差，即统计上所谓的“无偏估计”

## 分母为n造成的偏差

$$\begin{aligned}\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 &= \frac{1}{n} \sum_{i=1}^n \left[ (X_i - \mu) + (\mu - \bar{X}) \right]^2 \\&= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + \frac{2}{n} \sum_{i=1}^n (X_i - \mu)(\mu - \bar{X}) + \frac{1}{n} \sum_{i=1}^n (\mu - \bar{X})^2 \\&= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 + 2(\bar{X} - \mu)(\mu - \bar{X}) + (\mu - \bar{X})^2 \\&= \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 - (\mu - \bar{X})^2\end{aligned}$$

换言之, 除非正好  $\bar{X} = \mu$ , 否则我们一定有

$$\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 < \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2,$$

而不等式右边的那位才是对方的“正确”估计!

这个不等式说明了, 为什么直接使用  $\frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2$  会导致对方差的低估。

那么, 在不知道随机变量真实数学期望的前提下, 如何“正确”的估计方差呢? 答案是把上式中的分母  $n$  换成  $n - 1$ , 通过这种方法把原来的偏小的估计“放大”一点点, 我们就能获得对方差的正确估计了:

$$\mathbb{E} \left[ \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 \right] = \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \right] = \sigma^2.$$

$$S_1^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2.$$

但, 這個 estimator 有 bias, 因為:

$$\begin{aligned} E(S_1^2) &= \frac{1}{n} \sum_{i=1}^n E((X_i - \bar{X})^2) = \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu + \mu - \bar{X})^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n ((X_i - \mu)^2 - 2(X_i - \mu)(\bar{X} - \mu) + (\bar{X} - \mu)^2)\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - 2 \sum_{i=1}^n (X_i - \mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - 2n(\bar{X} - \mu)(\bar{X} - \mu) + n(\bar{X} - \mu)^2\right) \\ &= \frac{1}{n} E\left(\sum_{i=1}^n (X_i - \mu)^2 - n(\bar{X} - \mu)^2\right) \\ &= \frac{1}{n} \left(\sum_{i=1}^n E((X_i - \mu)^2) - nE((\bar{X} - \mu)^2)\right) \\ &= \frac{1}{n} (n\text{Var}(X) - n\text{Var}(\bar{X})) \\ &= \text{Var}(X) - \text{Var}(\bar{X}) = \sigma^2 - \frac{\sigma^2}{n} = \frac{n-1}{n} \sigma^2, \end{aligned}$$

而  $(n-1)/n * \sigma^2 \neq \sigma^2$ , 所以, 為了避免使用有 bias 的 estimator, 我們通常使用它的修正值  $S^2$ :

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

# 方差的Python例子

---

```
import numpy as np

arr = [1,2,3,4,5,6]
                                #求均值
arr_mean = np.mean(arr)
                                #求方差
arr_var = np.var(arr)
                                #求标准差
arr_std = np.std(arr,ddof=1)

print("平均值为: %f" % arr_mean)
print("方差为: %f" % arr_var)
print("标准差为:%f" % arr_std)
```

## 3.8 协方差

---

协方差用来刻画两个随机变量(X, Y)之间的相关性。

$$Cov(X, Y) = E[(X - E[X])(Y - E[Y])] = E[XY] - E[X]E[Y]$$

如果协方差为正，说明X, Y同向变化，协方差越大说明同向程度越高；如果协方差为负，说明X, Y反向运动，协方差越小说明反向程度越高。

# 协方差的python的例子

```
1. import numpy as np
2.
3. # 随机生成两个样本
4. x = np.random.randint(0, 9, 1000)
5. y = np.random.randint(0, 9, 1000)
6.
7. # 计算平均值
8. mx = x.mean()
9. my = y.mean()
10.
11. # 计算标准差
12. stdx = x.std()
13. stdy = y.std()
14.
15. # 计算协方差矩阵
16. covxy = np.cov(x, y)
17. print(covxy)
```

```
1. [[6.83907508 0.10925926]
2.  [0.10925926 6.53390891]]
3. 6.832236
4. 6.527375
5. 0.109149999999999989
6. [[1.          0.01634455]
7.  [0.01634455 1.          ]]
```