

Statistical Data Analysis

problem sheet 2



sat sun mon tue wed thu fri

Date: 04 / 11 / 2022

① (i) Given that,

$$f(t) = \begin{cases} \frac{1}{2\theta\sqrt{t}} \exp\left(-\frac{\sqrt{t}}{\theta}\right), & \text{for } t > 0 \\ 0 & \text{for } t \leq 0 \end{cases}$$

Now we call $f(t_1, \dots, t_n | \theta)$ as the likelihood function. It always depends on unknown parameter θ and it is denoted as $L(\theta)$.

$$\text{So, } L(\theta) = \prod_{i=1}^n f(t_i | \theta) \dots \dots \textcircled{i}$$

$L(\theta)$ is defined as a product of n terms, which is not easy to be maximized. So, maximizing $\log L(\theta)$ because \log is a monotonic increasing function. We define $\log L(\theta)$ as "log likelihood function",

$$l(\theta) = \log L(\theta) = \log \prod_{i=1}^n f(t_i | \theta)$$

$$= \sum_{i=1}^n \log f(t_i | \theta) \dots \dots \textcircled{ii}$$

Maximizing $l(\theta)$ with respect to θ will give us the MLE estimation.



sat sun mon tue wed thu fri

Date: / /

So, the log likelihood function,

$$l(\theta) = \log L(\theta) = \sum_{i=1}^n \log f(t_i | \theta) \quad [\text{using (ii)}]$$

$$= \sum_{i=1}^n \log \left\{ \frac{1}{2\theta\sqrt{t}} \exp\left(\frac{-\sqrt{t}}{\theta}\right) \right\}$$

$$= \sum_{i=1}^n \left[\log\left(\frac{1}{2\theta\sqrt{t}}\right) + \log\left(\exp\left(\frac{-\sqrt{t}}{\theta}\right)\right) \right]$$

$$= \sum_{i=1}^n \left[-\log 2\theta\sqrt{t} - \frac{\sqrt{t}}{\theta} \right]$$

$$= -\sum_{i=1}^n \log 2\theta\sqrt{t} - \sum_{i=1}^n \frac{\sqrt{t}}{\theta}$$

now, $l'(\theta) = -\sum_{i=1}^n \left\{ \frac{d}{d\theta} (\log 2\theta\sqrt{t}) - \frac{d}{d\theta} \left(\frac{\sqrt{t}}{\theta} \right) \right\}$

$$= -\sum_{i=1}^n \frac{1}{\theta} + \sum_{i=1}^n \frac{\sqrt{t}}{\theta^2}$$

$$\left[\frac{d}{dx} (\log x) = \frac{1}{x} \right]$$

$$\Rightarrow 0 = -\frac{n}{\theta} + \frac{1}{\theta^2} \sum_{i=1}^n \sqrt{t} \quad [\text{as derivative in } \theta]$$

$$\Rightarrow 0 = \frac{-n\theta + \sum_{i=1}^n \sqrt{t}}{\theta^2}$$

sat
○sun
○mon
○tue
○wed
○thu
○fri
○

Date:

/ /

$$\Rightarrow n\theta = \sum_{i=1}^n \sqrt{t_i}$$

$$\therefore \hat{\theta}_{MLE} = \frac{1}{n} \sum_{i=1}^n \sqrt{t_i}$$

This is the maximum likelihood estimate,

Now, for the given sample,

$$t_1 = 11300, t_2 = 5000, t_3 = 4300, t_4 = 8500$$

and $t_5 = 7900,$

$$\hat{\theta}_{MLE} = \frac{1}{5} \sqrt{11300 + 5000 + 4300 + 8500 + 7900}$$

$$= \frac{1}{5} (\sqrt{37000})$$

$$= 38.471$$



(ii) Method of moments is an estimation method of population parameters.

$$\text{Now, } \bar{T} = E(T) = 2\theta^2$$

$$\Rightarrow \theta^2 = \bar{T}/2$$

$$\therefore \theta = \sqrt{\bar{T}/2}$$

(2) Given that,

$$f_{\theta}(x) = \frac{\theta^x}{x!} e^{-\theta}, \quad x = 0, 1, 2, \dots$$

Maximum likelihood estimate,

$$\hat{\theta}_{MLE} = E\left(\frac{1}{n} \sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} E\left(\sum_{i=1}^n x_i\right)$$

$$= \frac{1}{n} \sum_{i=1}^n (E(x_i))$$

$$= \frac{1}{n} \sum_{i=1}^n \theta \quad [\because \text{it is poisson distribution function, } E(x) = \theta]$$

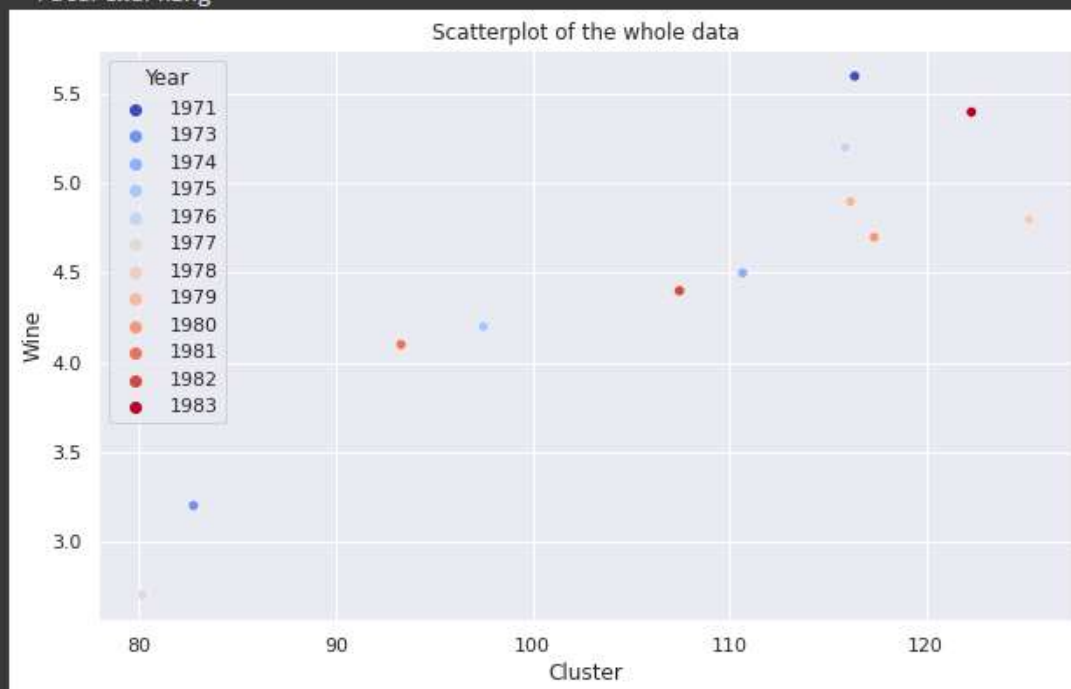
$$= \frac{1}{n} \cdot n\theta$$

$$= \theta$$

So, we see, that the estimate is unbiased.

```
[ ] 1 #setting style
2 sns.color_palette("Set2")
3 sns.set_style("whitegrid")
4 #define figure size
5 sns.set(rc={"figure.figsize":(10, 6)}) #width=8, height=4
6 #ploting scatterplot
7 sns.scatterplot(df["Cluster"],df["Wine"], df["Year"], legend="full", palette="coolwarm")
8 plt.title("Scatterplot of the whole data")
9
10 plt.show()
```

/usr/local/lib/python3.7/dist-packages/seaborn/_decorators.py:43: FutureWarning: Pass the following variables as keyword arguments: {'Year': 1971, 'Year': 1973, 'Year': 1974, 'Year': 1975, 'Year': 1976, 'Year': 1977, 'Year': 1978, 'Year': 1979, 'Year': 1980, 'Year': 1981, 'Year': 1982, 'Year': 1983}. This will ensure compatibility in future versions of Seaborn.



2. Assume the a simple linear regression model and estimate the parameters β_0 and β_1 in the regression model (use your own derivation).

```
[ ] 1 #defining x and y for the regression model
    2 x = df["Cluster"]
    3 y = df["Wine"]
    4
    5 N = len(x)
    6 #deriving a regression function
    7 def lin_reg(x,y):
    8     x_mean=x.mean()
    9     y_mean=y.mean()
10     x_var = x - x_mean
11     y_var = y - y_mean
12     #sum of products of x and y variance
13     sumof_xy_var = (x_var * y_var).sum()
14
15     sumof_sig_x = ((x_var)**2).sum()
16     #now we can estimate b1
17     b1 = sumof_xy_var/sumof_sig_x
18     #estimating b0
19     b0 = y_mean - (b1*x_mean)
20     #printing the line as y = mx + b = b1(x) + b0
21     simple_reg_lin = 'y = {} $\beta$  + {}'.format(round(b1, 4),round(b0, 4))
22
23     return(b0, b1, simple_reg_lin)
24
```

```
[ ] 1 #plugging in the data into the model
    2 lin_reg(x,y)
```

```
(-1.027902709556801, 0.051380577893637314, 'y = 0.0514 $\beta$  + -1.0279')
```

3. Plot the regression line

To find the regression line we need

- Correlation Coefficient (R)
- Coefficient of Determination (R^2)

We will be using pearson's correlation coefficient here to get the R and then simply square it to get the coefficient of determination which will tell us the goodness of fit.

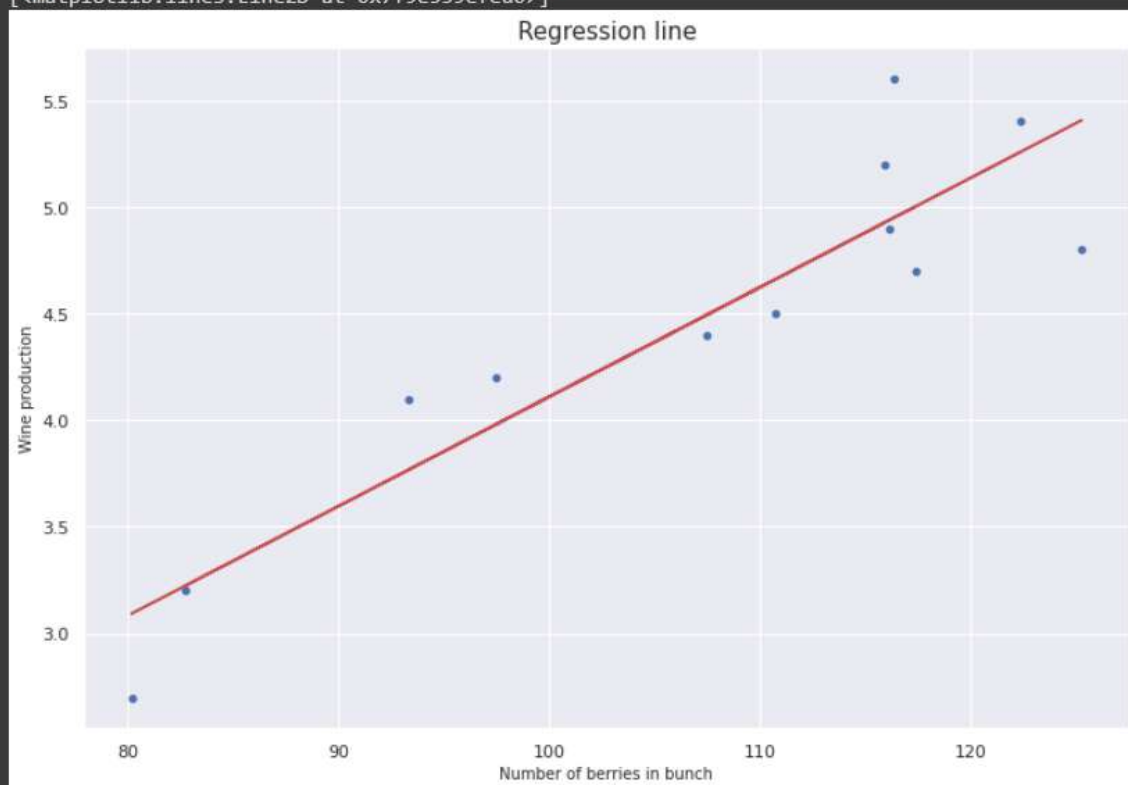
```
[ ] 1 #Pearson's correlation coefficient
2 def Pear_cor_coef(x,y):
3     #sum of the products of x and y
4     sumof_xy = (x * y).sum()
5     x_sum = x.sum()
6     y_sum = y.sum()
7     #product of summation of x and y
8     prod_sum = x_sum*y_sum
9     #getting the summation of x and y square
10    sumof_x_sq = (x**2).sum()
11    sumof_y_sq = (y**2).sum()
12    #now estimating R
13    r = ((N*sumof_xy)-((x_sum * y_sum))) / (np.sqrt(((N * sumof_x_sq)-(x_sum**2))*((N * sumof_y_sq)-(y_sum**2))))
14    return r
15
16 R = Pear_cor_coef(x,y)
17 #now estimating coefficient of determination
18 coef_det = R**2
19
```

```
[ ] 1 #printing the correaltion coefficient and coefficient of determination
2 print('Correlation Coefficient: ', R)
3 print('Coefficient of Determination: ', coef_det)
```

```
Correlation Coefficient:  0.913214127416157
Coefficient of Determination:  0.8339600425124531
```

```
[ ] 1 #now plotting the regression line
2 plt.figure(figsize=(12,8))
3 sns.scatterplot(x, y, palette="inferno")
4 plt.title('Regression line', fontsize=15)
5 plt.xlabel('Number of berries in bunch', fontsize=10)
6 plt.ylabel('Wine production', fontsize=10)
7 b0, b1, regline = lin_reg(x,y)
8 plt.plot(x, b0 + b1*x, c='r')
```


[<matplotlib.lines.Line2D at 0x7f9e539efed0>]



4. Predict the yearly production of wine if the number of berries in a bunch of grapes is 100.

```
[ ] 1 #predicting the value of wine production when the number of berries in a bunch of grapes is 100
    2 newy = b0 + (b1*100)
    3 print("Predicted Wine Production : ", round(newy, 2))
```

```
Predicted Wine Production : 4.11
```

So The Predicted yearly wine production is 4.11 tons per 100 meter square if the number of berries in a bunch of grapes is 100.