

ass1

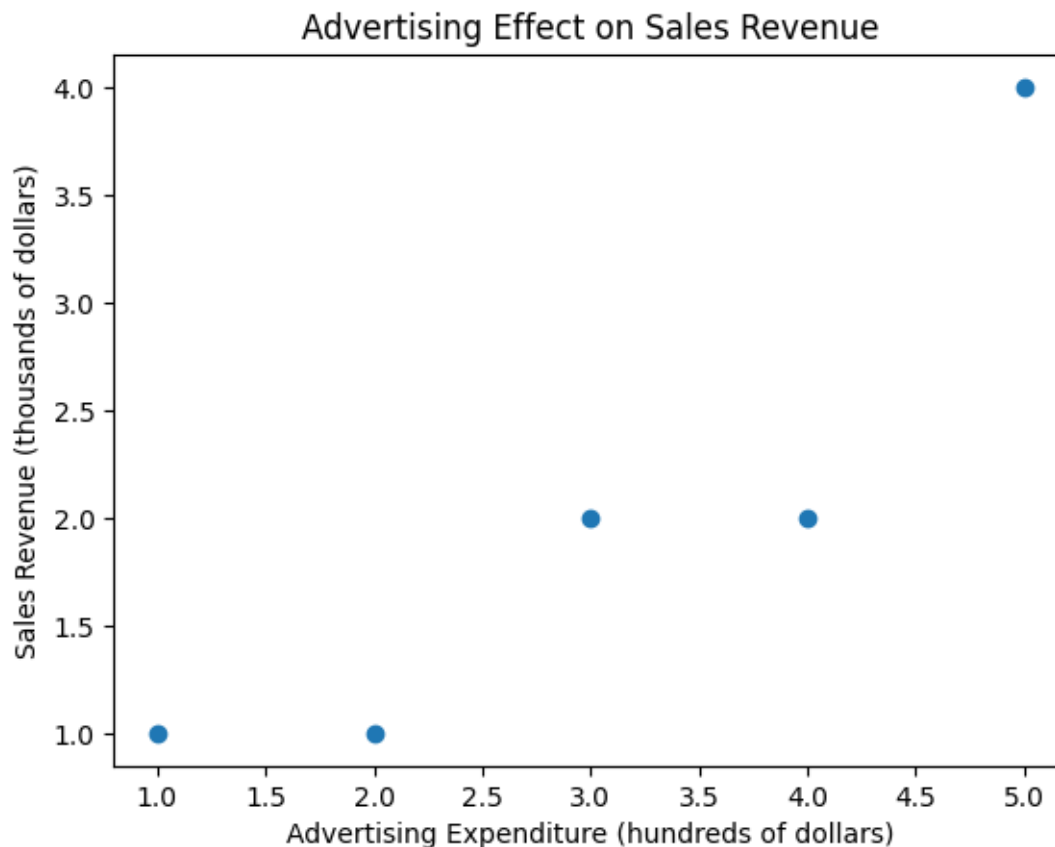
September 23, 2024

Suppose an appliance store conducts a 5-month experiment to determine the effect of advertising on sales revenue. The results are shown below. Advertising Expenditure x (hundreds of dollars) 1 2 3 4 5 Sales Revenue y (thousands of dollars) 1 1 2 2 4

(a) Draw a scatterplot of the data and comment the relationship between y and x .

```
[30]: import matplotlib.pyplot as plt
```

```
[31]: x=[1,2,3,4,5]
      y=[1,1,2,2,4]
      plt.scatter(x,y)
      plt.xlabel('Advertising Expenditure (hundreds of dollars)')
      plt.ylabel('Sales Revenue (thousands of dollars)')
      plt.title('Advertising Effect on Sales Revenue')
      plt.show()
```



as we can see from the scatter plot, there is a positive linear relationship between advertising expenditure and sales revenue.

(b) What is your linear regression model? State the necessary assumptions.

The linear regression model is:

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

where β_0 is the intercept and β_1 is the slope.

The necessary assumptions for linear regression are:

1. $\epsilon_i \sim \mathcal{N}(\mu, \sigma^2)$
2. The errors are independent of each other.

(c) Find the least squares line from the data and plot it on your scatterplot.

we can calculate the parameters using the following formulas: $\beta_1 = \frac{n \sum (x_i y_i) - \sum x_i \sum y_i}{n \sum (x_i^2) - (\sum x_i)^2}$ $\beta_0 = \bar{y} - \beta_1 \bar{x}$

```
[32]: from scipy.stats import linregress
import numpy as np
```

```
[33]: def sum_of_array(x):  
      ans = 0  
      for i in x:  
          ans += i  
      return ans
```

```
[34]: def sum_of_squares(x):  
      ans = 0  
      for i in x:  
          ans += i**2  
      return ans
```

```
[35]: def sum_of_xy( x , y):  
      ans = 0  
      for i in range(len(x)):  
          ans += x[i] * y[i]  
      return ans
```

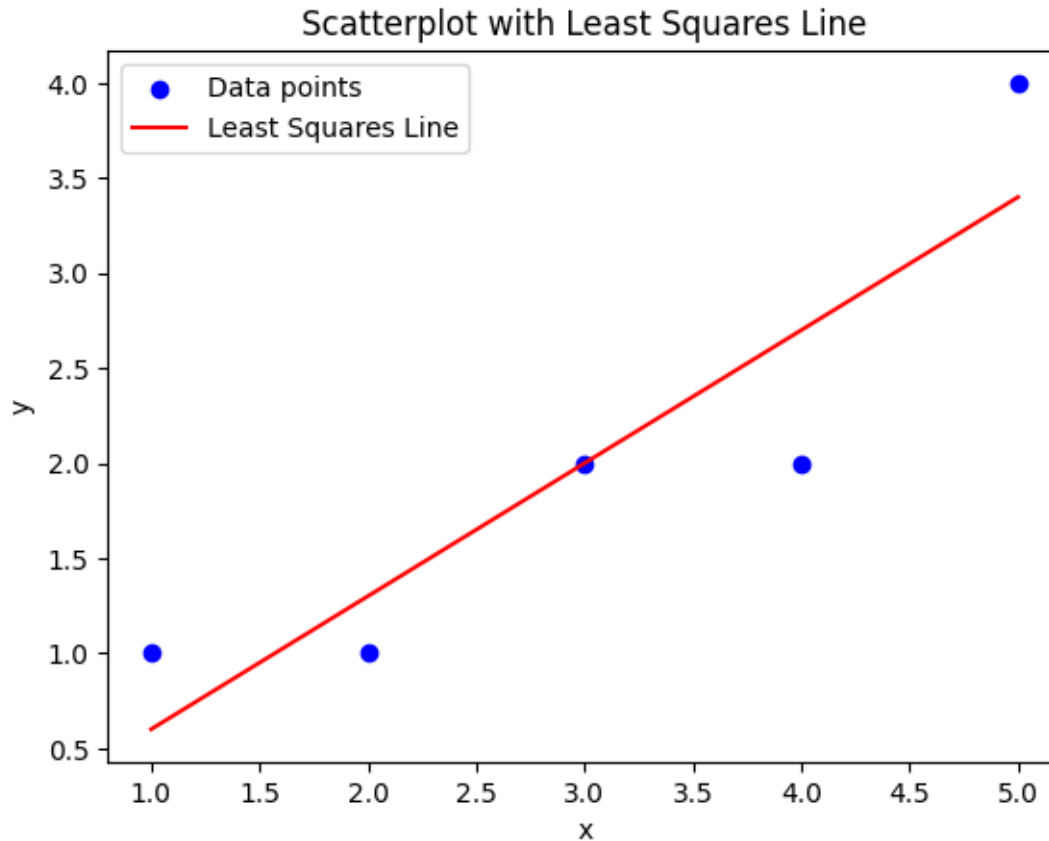
```
[36]: Sxy=(len(x) * sum_of_xy(x, y) - sum_of_array(x) * sum_of_array(y))  
      Sxx=(len(x) * sum_of_squares(x) - sum_of_array(x)**2)  
      def cal_beta1():  
          ans= Sxy/ Sxx  
          return ans
```

```
[37]: def cal_beta0(x, y):  
      ans = sum_of_array(y) / len(x) - cal_beta1() * sum_of_array(x) / len(x)  
      # print(f"beta0 = { ans}")  
      return ans
```

```
[38]: slope = cal_beta1()  
      intercept = cal_beta0(x, y)  
      print('slope:', slope)  
      print('intercept:', intercept)  
      plt.scatter(x, y, color='blue', label='Data points')  
      regression_line = slope * np.array(x) + intercept  
      plt.plot(x, regression_line, color='red', label='Least Squares Line')  
      plt.xlabel('x')  
      plt.ylabel('y')  
      plt.title('Scatterplot with Least Squares Line')  
      plt.legend()  
      plt.show()
```

slope: 0.7

intercept: -0.10000000000000009



(d) Test the hypothesis that the Advertising Expenditure has no effect of the Sales Revenue when a linear model is used (use $\alpha = 0.05$). State the null and alternative hypotheses. Draw the appropriate test conclusions.

for β_1 : $H_0 : \beta_1 = 0$ $H_1 : \beta_1 \neq 0$ Test statistic: $t = \frac{\hat{\beta}_1 - 0}{S/\sqrt{S_{xx}}}$ $t \sim t_{n-2}$ if H_0 is true $\text{Decision rule: reject } H_0 \text{ if } |t| > t_{\alpha/2, n-2}$.

```
[39]: t_value = (slope - 0) / (np.sqrt(Sxy / Sxx)) # t(0.05, 4) = 2.7764
      print('t-value:', t_value)
```

t-value: 0.8366600265340755

because $t_{\alpha/2, n-2} = 2.7764$, t-value: 0.8366600265340755, $|t| < t_{\alpha/2, n-2}$ we not reject H_0

(e) Find a 95% confidence interval for β_1 (slope of the linear regression model). Interpret your results.

the 95% confidence interval for β_1 is:

$$\hat{\beta}_1 \pm t_{\alpha/2, n-2} S/\sqrt{S_{xx}}$$

where $t_{\alpha/2, n-2}$ is the t-value corresponding to the 95% confidence level.

$$\hat{\beta}_1 = 0.7$$

The estimator of σ^2 is given by:

$$S^2 = \text{MSE} = \text{SSE} \frac{1}{n-2} = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n-2}$$

Here, S^2 is an unbiased estimator of σ^2 .

```
[40]: y_ave= np.average(y)/len(y)
y_eva=[intercept + slope*x[i] for i in range(len(x))]
S=sum((y[i]-y_eva[i])**2 for i in range(len(x)))/(len(x)-2)
CI=(slope - 2.7764 *S/ np.sqrt(Sxx), slope + 2.7764*S / np.sqrt(Sxx))
print('95% confidence interval for 1:', CI)
```

95% confidence interval for 1: (0.5560311737323358, 0.8439688262676641)

because $\beta_1 = 0.7 \sim (0.5560311737323358, 0.8439688262676641)$ so we can say that the effect of advertising expenditure on sales revenue is significant at 95% confidence level.

(f) Find the coefficient of determination for the linear regression model. Interpret your result.

The coefficient of determination, or R^2 , is defined as $R^2 = \text{SSR} \frac{1}{\text{SST} - \frac{\text{SSE}}{\text{SST}} \text{SST}} = \frac{\text{SSR}}{\text{SST}}$

\$

where

\$

$$\begin{aligned} \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_0} &= \sum_{i=1}^n \frac{y_i - \beta_0 - \beta_1 x_i}{\sigma^2}, \\ \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \beta_1} &= \sum_{i=1}^n \frac{(y_i - \beta_0 - \beta_1 x_i) x_i}{\sigma^2}, \\ \frac{\partial \ell(\beta_0, \beta_1, \sigma^2)}{\partial \sigma^2} &= \frac{1}{2\sigma^4} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 - \frac{n}{2\sigma^2} \end{aligned}$$

\$

Let

$$\ell(\beta_0, \beta_1, \sigma^2) \Big|_{\beta_0=0}$$

and get

$$\ell_0 = \bar{y} - \beta_1 \bar{x}$$

so that

$$\ell(\beta_0, \beta_1, \sigma^2) \Big|_{\beta_1=\sum_{i=1}^n \frac{x_i(y_i - \bar{y}) + \beta_1 x_i(\bar{x} - x_i)}{\sigma^2}}$$

Let

$$\ell(\beta_0, \beta_1, \sigma^2) \Big|_{\beta_1=0}$$

and the MLE of β_1 is

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i(y_i - \bar{y})}{\sum_{i=1}^n x_i(x_i - \bar{x})} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

the MLE of β_0 is

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} = \hat{y} - \frac{\sum_{i=1}^n x_i(x_i - y)}{\sum_{i=1}^n x_i(x_i - \bar{x})}$$

Let

$$l(\beta_0, \beta_1, \sigma^2) \big|_{\partial \sigma^2 = 0}$$

and get the MLE of σ^2 is

$$\hat{\sigma}^2 = \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{n}$$

where

$$\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$$

(b) Comments on the difference between MLE and LSE (least square estimation)

The main difference between Maximum Likelihood Estimation (MLE) and Least Squares Estimation (LSE) is that MLE estimates parameters by maximizing the likelihood function to make the observed data most probable, whereas LSE estimates parameters by minimizing the sum of squared differences between observed and predicted values, typically used in linear regression; additionally, MLE is generally sensitive to model assumptions, while LSE can perform poorly in the presence of outliers.

[42] :