

$$\text{Corr}(g_1, g_2) = \frac{\text{cov}(g_1, g_2)}{\sigma(g_1) \sigma(g_2)}$$

$$\bar{g}_1 = 0 \rightarrow g_1 - \bar{g}_1 = (-5 \ -4 \ -3 \ -2 \ -1 \ 0 \ 1 \ 2 \ 3 \ 4 \ 5) \quad (1)$$

$$\bar{g}_2 = 11 \rightarrow g_2 - \bar{g}_2 = (14 \ 5 \ -2 \ -7 \ -10 \ -10 \ -7 \ -2 \ 5 \ 14) \quad (2)$$

$$\text{cov}(g_1, g_2) = \frac{1}{n-1} \sum_{k=1}^n (g_{1k} - \bar{g}_1)(g_{2k} - \bar{g}_2)$$

$$= \frac{1}{9} (14 \times -5 + 5 \times -4 + \dots + 5 \times 4 + 14 \times 5) = 0$$

$$\rightarrow \text{Corr}(g_1, g_2) = \frac{\text{cov}(g_1, g_2)}{\sigma(g_1) \sigma(g_2)} = \frac{0}{\dots} = 0$$

پس با استفاده از این مقیار این ۲ داده هیچ شباهتی بهم ندارند و در نتیجه قرار نمی گیرند.

Mutual Information (b)

$$I(g_1, g_2) = H(g_1) + H(g_2) - H(g_1, g_2)$$

$$H(g_1) = - \sum_{i=1}^{10} P_i \log P_i = - \sum \frac{1}{10} \log \frac{1}{10} = - \log \frac{1}{10} = \log 10$$

$$H(g_2) = - \sum_{i=1}^5 P_i \log P_i = - \sum \frac{1}{5} \log \frac{1}{5} = - \log \frac{1}{5} = \log 5$$

$$H(g_1, g_2) = - \sum_{i,j} P_{ij} \log P_{ij} = 10 \text{ مورد } 10 \text{ برابر } 10 \text{ و برابر } 10$$

$$= -10 \left(\frac{1}{10} \log \frac{1}{10} \right) = - \log \frac{1}{10} = \log 10$$

$$I(g_1, g_2) = \log 10 + \log 5 - \log 10 = \log 5$$

ج) بله از آن جا که داده ها رابطه ی خطی ندارند بنابراین Corr آن ۲ صفر شده چون Corr رابطه ی خطی را نشان می دهد (اگر نه)

اما $M2$ رابطه ی غیر خطی را نیز نشان می دهد و ۵ زده است

پس اگر g_2 رابطه ی غیر خطی با g_1 داشته باشد

a → euclidean: $\sqrt{(2-1)^2 + 1^2 + 1^2 + 1^2} = \sqrt{4} = \boxed{2}$

(2)

→ correlation: $\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \boxed{0}$

$$\text{cov}(x, y) = \frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})(y_k - \bar{y})$$

$$= \frac{1}{4} (0 + 0 + 0 + 0) = 0$$

→ cosine: $\cos(x, y) = \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{2+2+2+2}{\sqrt{4+4+4+4} \times \sqrt{4}} = \frac{8}{8 \times 2} = \boxed{1}$

b → jaccard: $J = \frac{f_{11}}{f_{01} + f_{10} + f_{11}} = \frac{0}{0+0+0} = \boxed{0}$

→ euclidean: $\sqrt{1^2 + 1^2 + 1^2 + 1^2} = \sqrt{4} = \boxed{2}$

→ correlation: $\text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$

$x - \bar{x} = (\frac{1}{4}, -\frac{1}{4}, \frac{1}{4}, -\frac{1}{4})$

$y - \bar{y} = (-\frac{1}{4}, \frac{1}{4}, -\frac{1}{4}, \frac{1}{4})$

$$\text{cov}(x, y) = \frac{1}{4} (-\frac{1}{4} - \frac{1}{4} - \frac{1}{4} - \frac{1}{4}) = -\frac{1}{4}$$

$$\sigma_x = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (x_k - \bar{x})^2} = \sqrt{\frac{1}{4} (\frac{1}{4} \times 4)} = \sqrt{\frac{1}{4}}$$

$$\sigma_y = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (y_k - \bar{y})^2} = \sqrt{\frac{1}{4}}$$

$$\text{corr}(x, y) = \frac{-\frac{1}{4}}{\sqrt{\frac{1}{4}} \sqrt{\frac{1}{4}}} = \boxed{-1}$$

→ cosine: $\cos(x, y) = \frac{\langle x, y \rangle}{\|x\| \|y\|} = \frac{0}{2 \times 2} = \boxed{0}$

c) \rightarrow bhattacharyya $(x, y) = D_B(x, y) = -\ln\left(\sum_i \sqrt{x_i y_i}\right)$
 $= -\ln(\sqrt{1} + \sqrt{1} + \sqrt{0} + \sqrt{0} + \sqrt{0} + \sqrt{1})$
 $= -\ln(3) = \boxed{}$

\rightarrow Correlation: $\text{Corr}(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$

$$\bar{x} = \frac{4}{10}$$

$$x - \bar{x} = \left(\frac{1}{10}, \frac{1}{10}, \frac{1}{10}, -\frac{2}{10}, -\frac{2}{10}, \frac{1}{10}\right)$$

$$\bar{y} = \frac{4}{10} \rightarrow y - \bar{y} = \left(\frac{1}{10}, \frac{1}{10}, \frac{1}{10}, \frac{1}{10}, -\frac{2}{10}, \frac{1}{10}\right)$$

$$\text{Cov}(x, y) = \frac{1}{10} \left(\frac{1}{10} + \frac{1}{10} - \frac{2}{10} - \frac{2}{10} + \frac{4}{10} + \frac{1}{10} \right) = \frac{1}{10}$$

$$\sigma_x = \sqrt{\frac{1}{10} \left(\frac{1}{10} + \frac{1}{10} + \frac{4}{10} + \frac{4}{10} + \frac{1}{10} \right)} = \sqrt{\frac{1}{10} \left(\frac{12}{10} \right)} = \sqrt{\frac{6}{10}}$$

$$\sigma_y = \sqrt{\frac{1}{10} \left(\frac{1}{10} + \frac{1}{10} + \frac{1}{10} + \frac{1}{10} + \frac{4}{10} + \frac{1}{10} \right)} = \sqrt{\frac{6}{10}}$$

$$\text{Corr}(x, y) = \frac{\frac{1}{10}}{\sqrt{\frac{6}{10}} \sqrt{\frac{6}{10}}} = \frac{\frac{1}{10}}{\frac{6}{10}} = \boxed{\frac{1}{6}}$$

\rightarrow manhattan = 1 pi

$$(0 + 0 + 1 + 1 + 0 + 0) = \boxed{2}$$

۳

(a) بیروستہ - کمی - نسبت

(b) گہستہ - کیفی - ترتیبی

(c) بیروستہ - کمی - نسبت

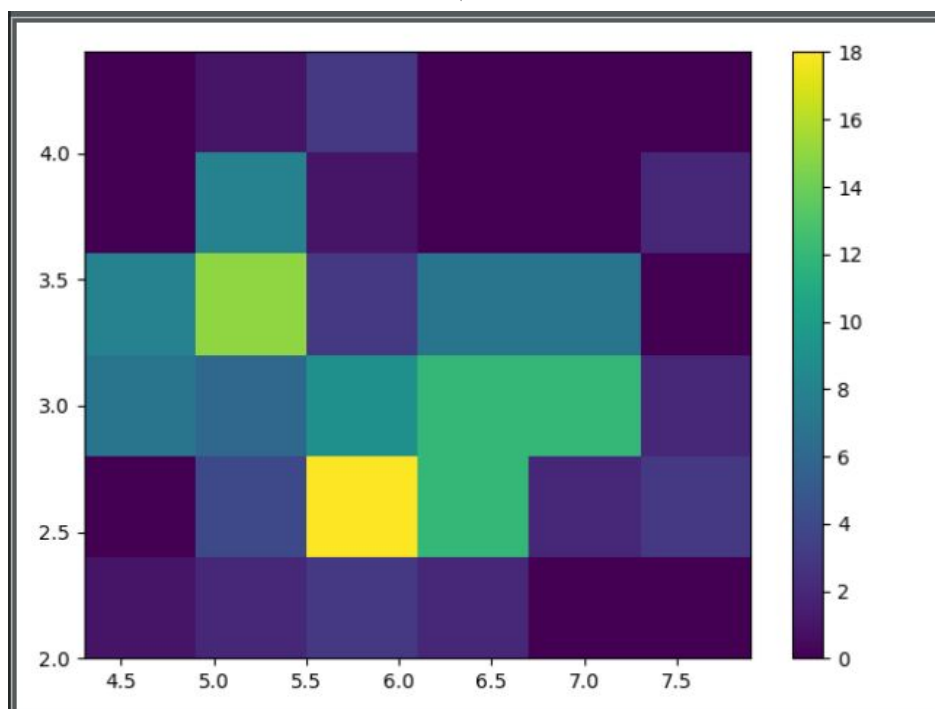
(d) بیروستہ - کمی - بازہ

(e) گہستہ - کیفی - ترتیبی

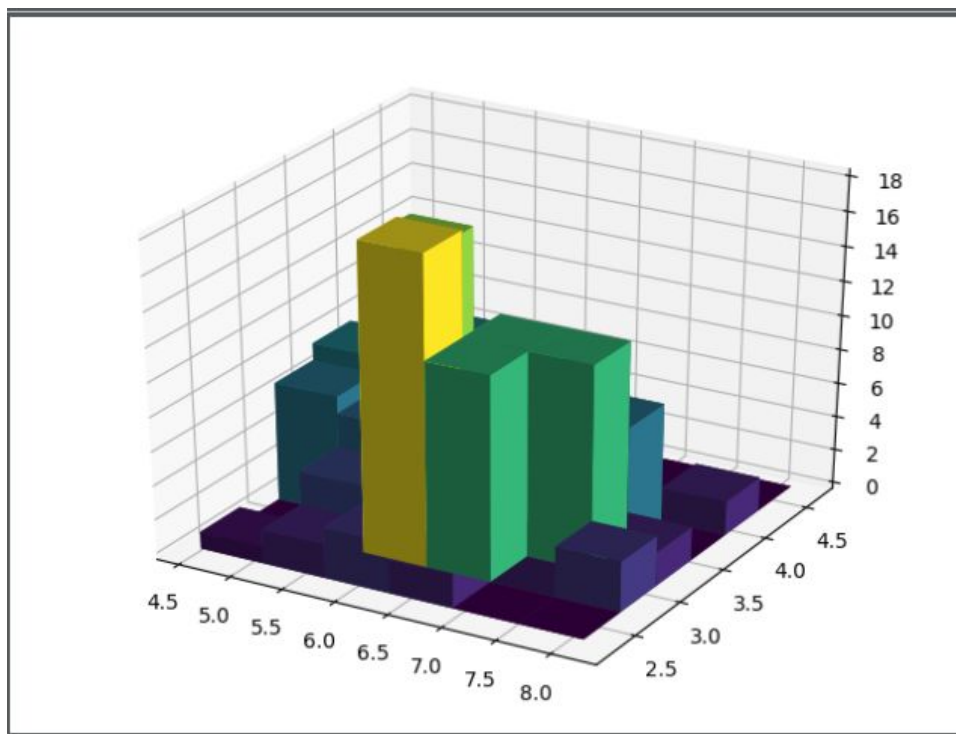
(f) آلتہ - کیفی - ترتیبی

جملہ ناظم انسانی تہا ~~تاریک~~ روشن و متبادل رات صیفی سے

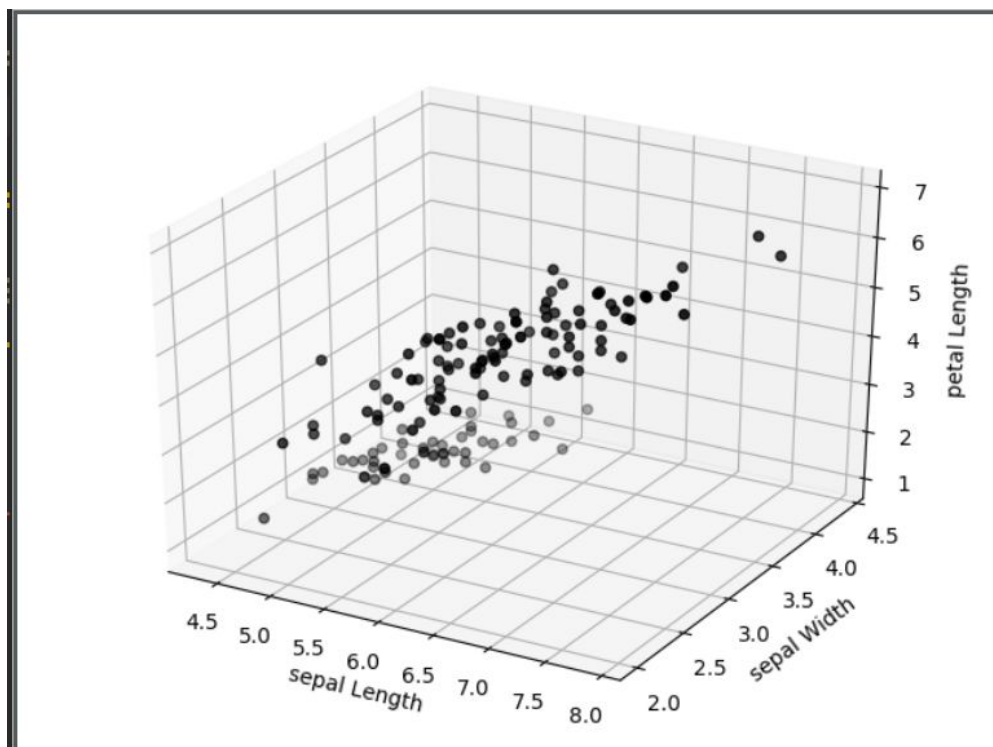
(a) دو ویژگی اول انتخاب شده اند. محور x بیانگر فیچر اول و محور y بیانگر فیچر دوم است. رنگ نزدیکتر به زرد بیانگر آن است که تعداد بیشتری داده دارای فیچر اول و فیچر دوم برابر با مقدار مشخص شده وجود دارند



(b) دو ویژگی اول انتخاب شده اند. محور x بیانگر فیچر اول و محور y بیانگر فیچر دوم است. محور عمودی هرچه بالاتر باشد بیانگر آن است که تعداد بیشتری داده دارای فیچر اول و فیچر دوم برابر با مقدار مشخص شده وجود دارند



(c) سه ویژگی اول انتخاب شده اند. محور x بیانگر فیچر اول و محور y بیانگر فیچر دوم است و محور z بیانگر فیچر سوم است و هر نقطه بیانگر یک داده است. هر محور با استفاده از label مشخص شده مربوط به کدام اتریبیوت است.



(d)

```
Mean : [5.843333333333334, 3.0540000000000003, 3.7586666666666666, 1.1986666666666668]
Variance : [0.6811222222222223, 0.18675066666666668, 3.092424888888889, 0.5785315555555555]
```

(e) دو فیچر اول انتخاب شده اند نتیجه یک آرایه ی ۲ در ۲ است که درایه ی i,j بیانگر کواریانس فیچر iam با jam است.

```
[[ 0.68569351 -0.03926846]
 [-0.03926846  0.18800403]]
```

(f) همانند e

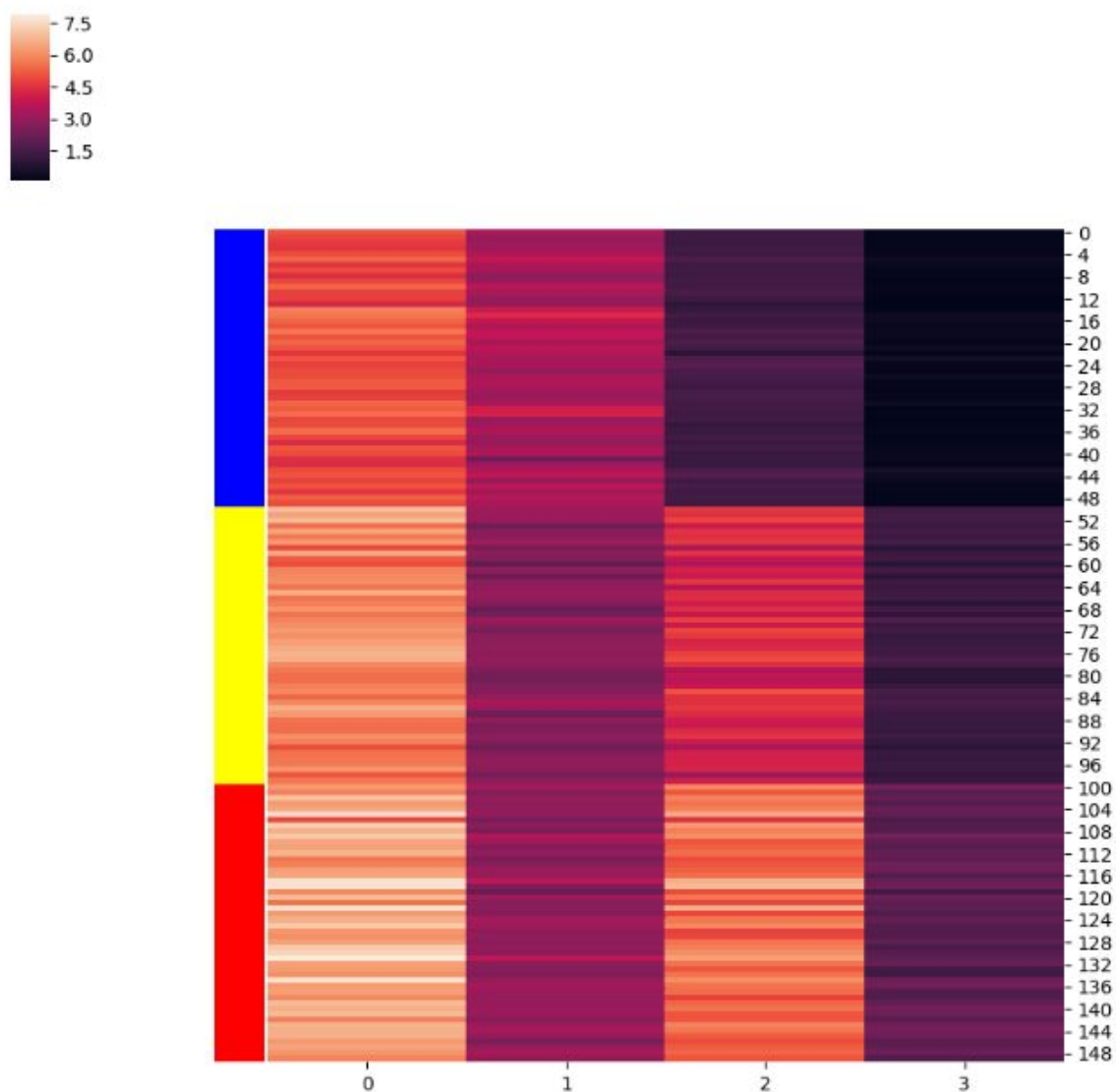
```
[[ 1. -0.10936925]
 [-0.10936925  1. ]]
```

(g) ماتریس همبستگی فیچر ها را برای گیاه مدنظر بدست می آوریم هر درایه که عدد بیشتری داشته باشد فیچر های آن وابستگی بیشتری دارند.

```
[[1.          0.45722782 0.86422473 0.28110771]
 [0.45722782 1.          0.40104458 0.53772803]
 [0.86422473 0.40104458 1.          0.32210822]
 [0.28110771 0.53772803 0.32210822 1.          ]]
```

فیچر های اول و سوم عدد بزرگتری دارد پس بیشترین شباهت برای این ۲ فیچر است

(h)

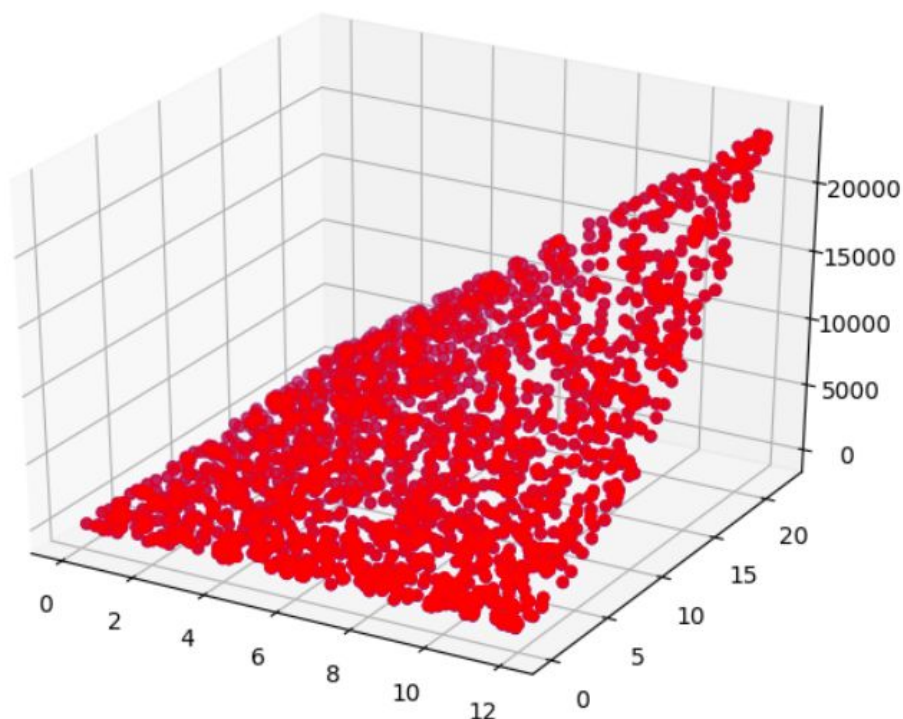


Clustermap کشیده شده نشان می‌دهد که فیچر سوم و چهارم بهتر از ۲ فیچر دیگر می‌توانند بین گیاهان مختلف تمایز قایل شوند و تفکیک کنند. همچنین به نظر می‌رسد فیچر سوم بهترین تفکیک گیاهان از یکدیگر را داشته باشد.

5- همانطور که در صورت سوال ذکر شده از مدل صورت سوال استفاده کرده ام و ضرایب آن را با **train** کردن بدست آورده ام که بسیار نزدیک به واقعیت شده است.
برای **gradient descent** :

```
Reading data ...  
Data read!  
Wait for training. It can take a few seconds!  
Beta = [[0.01608638]  
[0.0673921 ]  
[2.0030327 ]  
[4.00727882]]  
SSE = 1371227.3654922154  
SSE for train data : 1371227.3654922154  
SSE for test data : 425516.4013891396  
  
Process finished with exit code 0
```

مقایسه ی y واقعی داده ها و y پیش بینی شده (امتیازی):
 Y واقعی با رنگ قرمز و y پیش بینی شده با رنگ آبی مشخص شده. مشخص است که داده ها کاملاً روی هم قرار گرفته اند.



برای stochastic gradient descent :

```
Reading data ...  
Data read!  
Wait for training. It can take a few seconds!  
Beta = [[0.94145994]  
        [3.01608337]  
        [2.00105947]  
        [3.99976178]]  
SSE = 70.97660980508954  
SSE for train data : 70.97660980508954  
SSE for test data : 52.98509267995236  
  
Process finished with exit code 0
```

مقایسه ی y واقعی داده ها و y پیش بینی شده (امتیازی):
 Y واقعی با رنگ قرمز و y پیش بینی شده با رنگ آبی مشخص شده. مشخص است که داده ها کاملاً روی هم قرار گرفته اند.

