

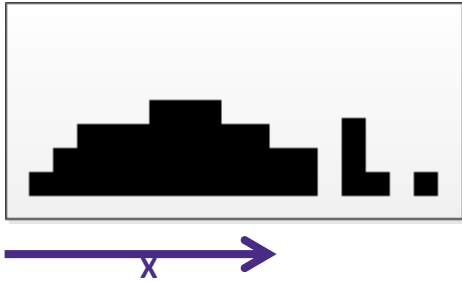
PROFESSIONAL & CONTINUING EDUCATION

UNIVERSITY *of* WASHINGTON

DIMENSIONS IN CLUSTERING



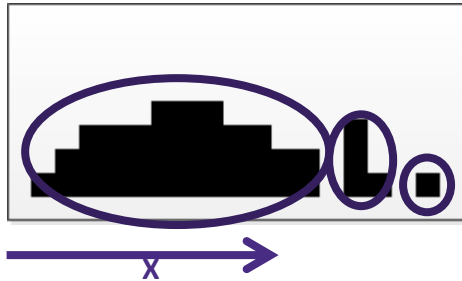
CLUSTERING: DIMENSIONS (1)



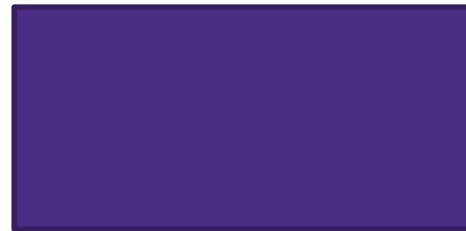
Where are the three clusters?

W

CLUSTERING: DIMENSIONS (2)

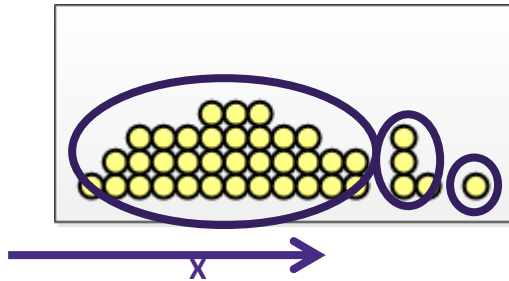


Simple assignment
based on a 1D
distribution



W

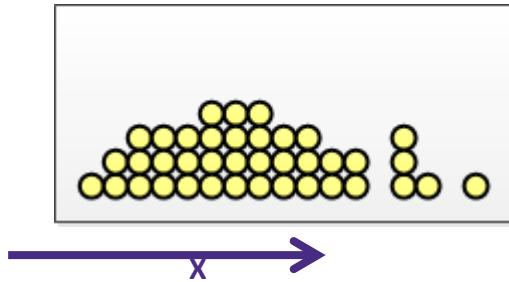
CLUSTERING: DIMENSIONS (3)



Simple assignment
based on a 1D
distribution

W

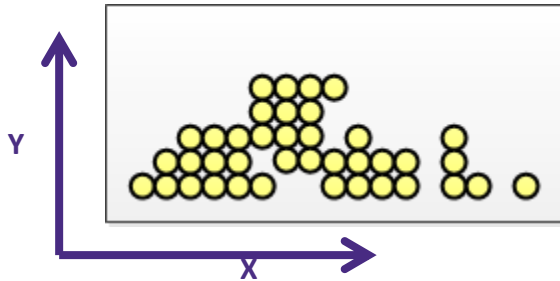
CLUSTERING: DIMENSIONS (4)



What if this was not
a 1D distribution?

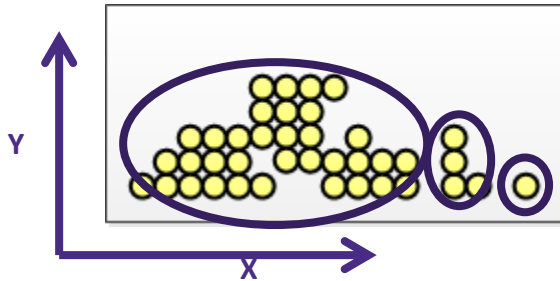
W

CLUSTERING: DIMENSIONS (5)



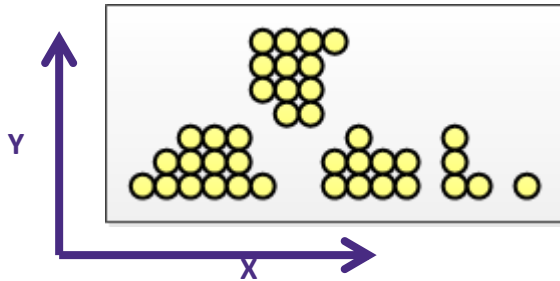
The distribution is in 2D. Some points differ in the 2nd D

CLUSTERING: DIMENSIONS (6)



If the difference is minor, we still get the same clusters

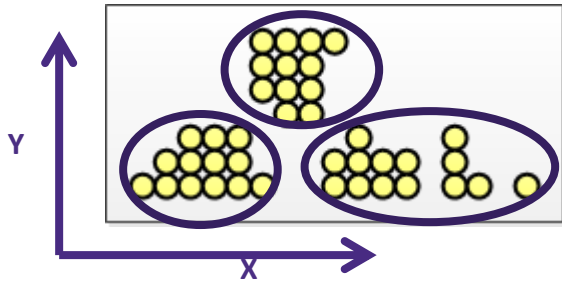
CLUSTERING: DIMENSIONS (7)



The difference could
be significant

W

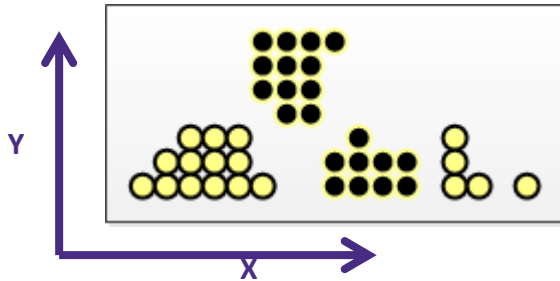
CLUSTERING: DIMENSIONS (8)



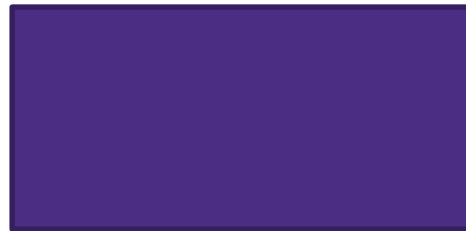
A big difference in the 2nd D can lead to different clusters

W

CLUSTERING: DIMENSIONS (9)

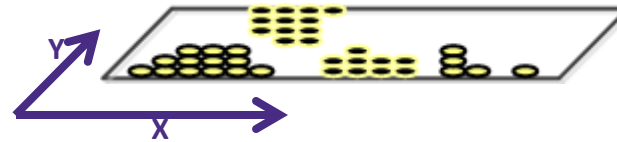
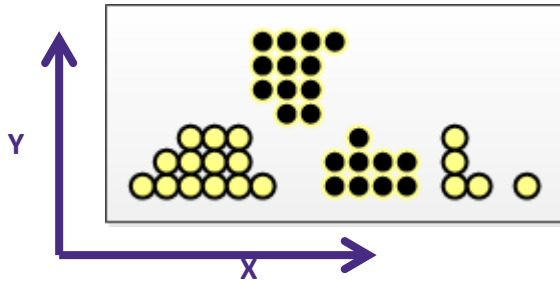


We can introduce another D by color coding. This is a Boolean Dimension



W

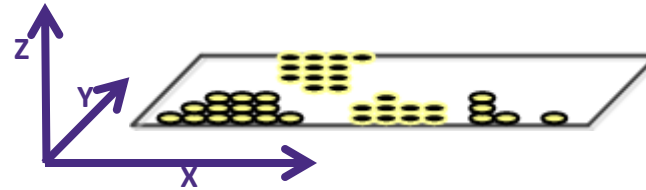
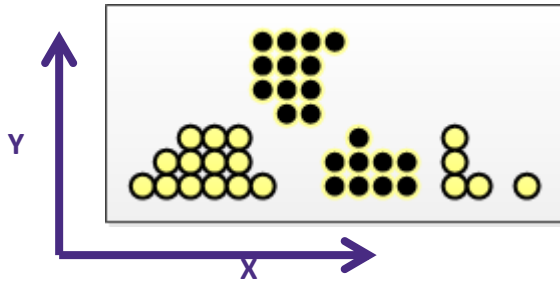
CLUSTERING: DIMENSIONS (10)



Create a 3rd
Dimension

W

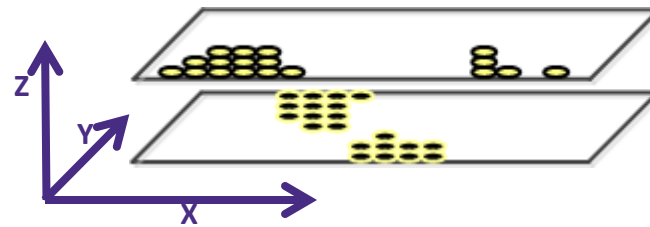
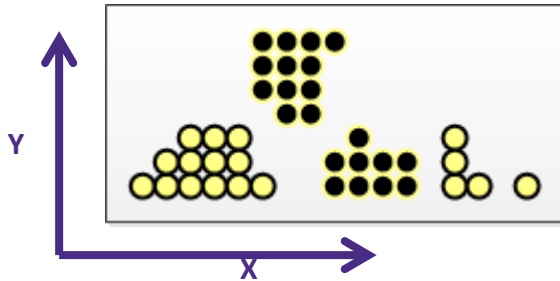
CLUSTERING: DIMENSIONS (11)



Create a 3rd
Dimension

W

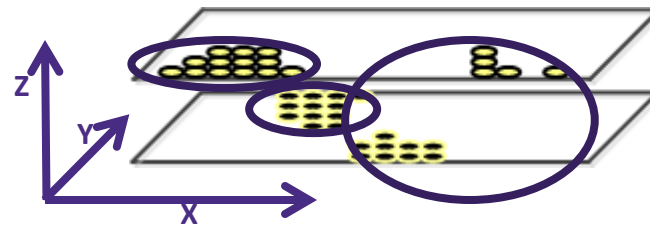
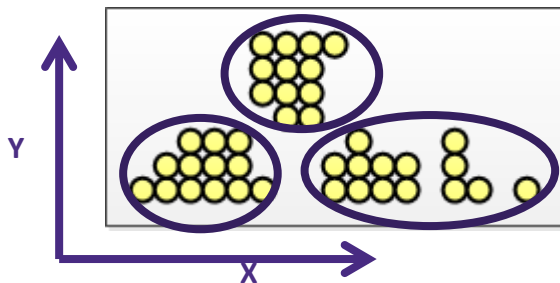
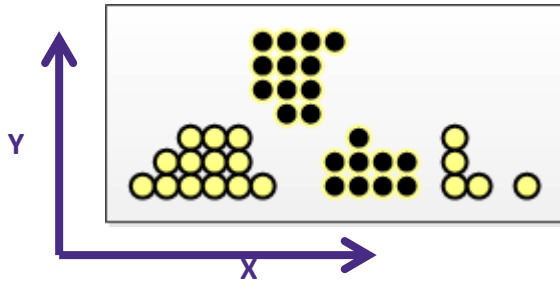
CLUSTERING: DIMENSIONS (12)



Where are the 3
clusters now?

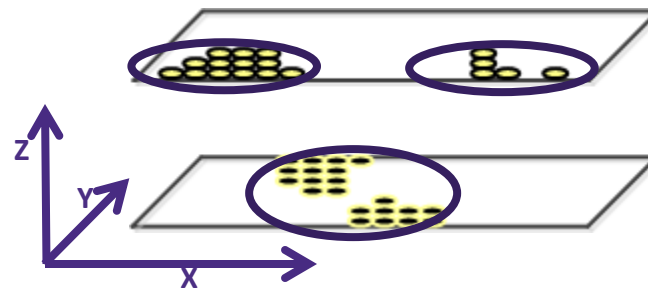
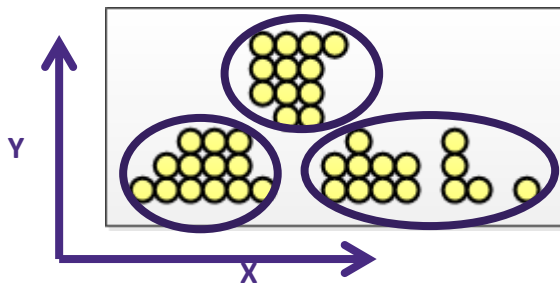
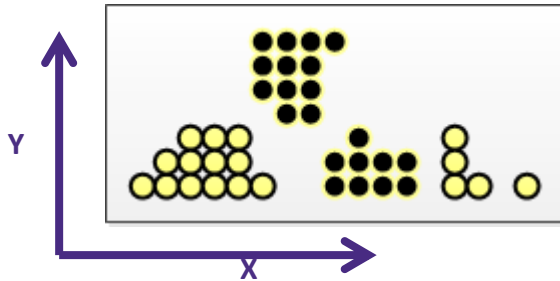
W

CLUSTERING: DIMENSIONS (13)



If the 3rd is small,
then the clustering is
the same as in 2D

CLUSTERING: DIMENSIONS (14)



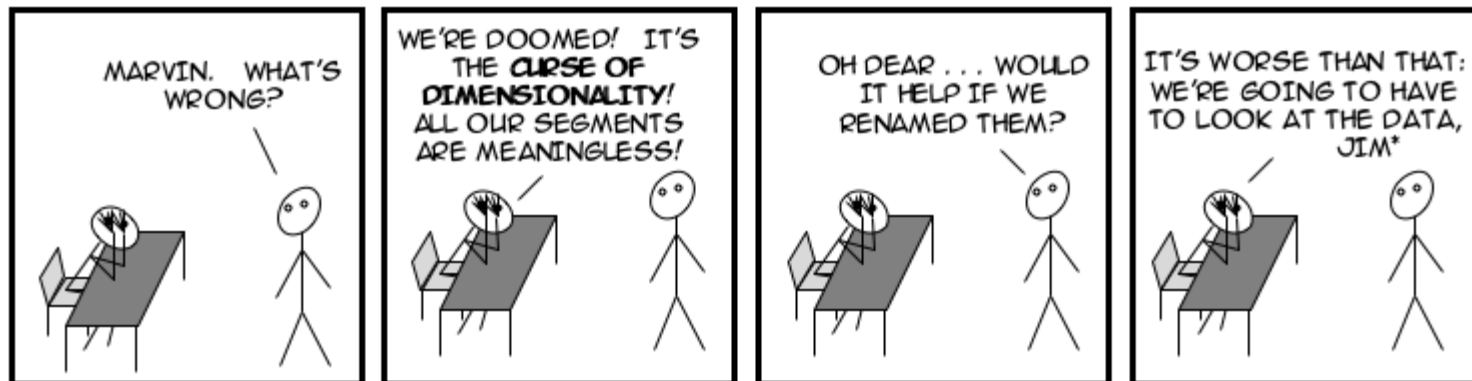
If the 3rd is big, then
the clustering differs
from 2D

W

Dimensions in Clustering



BREAK



[HTTP://SCIENTIFICMARKETER.COM](http://scientificmarketer.com)

COPYRIGHT © NICHOLAS J RADCLIFFE 2007. ALL RIGHTS RESERVED.
* WITH APOLOGIES TO MR SPOCK & STAR TREK.

W

Normalization in Clustering

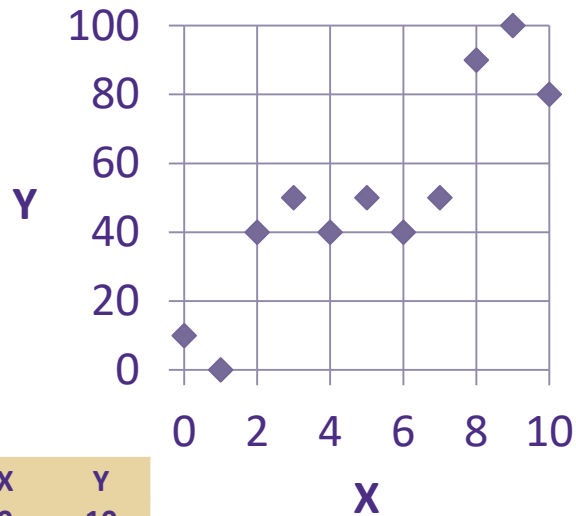


NORMALIZATION OF A LINEAR RELATIONSHIP (1)

X	Y
0	10
1	0
2	40
3	50
4	40
5	50
6	40
7	50
8	90
9	100
10	80



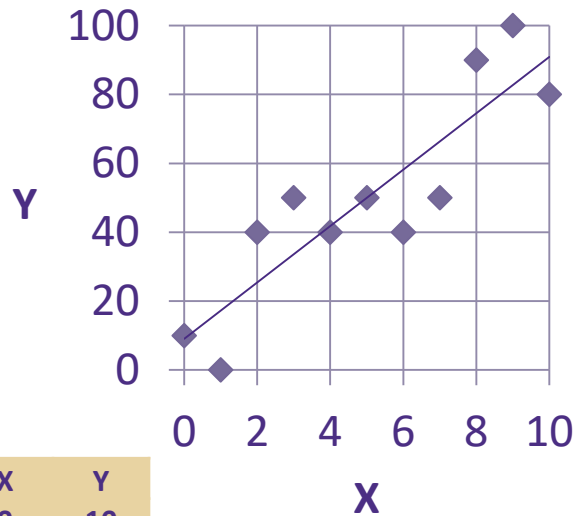
NORMALIZATION OF A LINEAR RELATIONSHIP (2)



X	Y
0	10
1	0
2	40
3	50
4	40
5	50
6	40
7	50
8	90
9	100
10	80



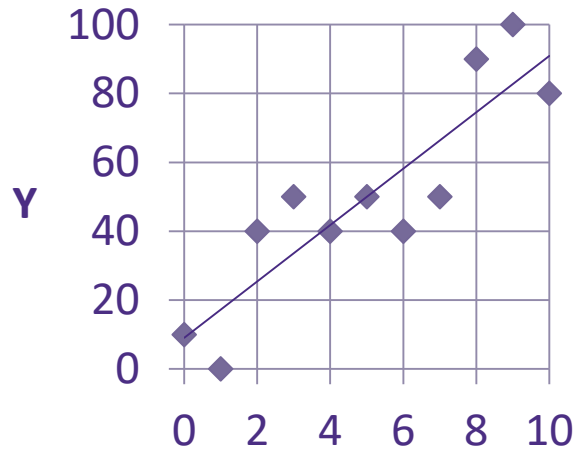
NORMALIZATION OF A LINEAR RELATIONSHIP (3)



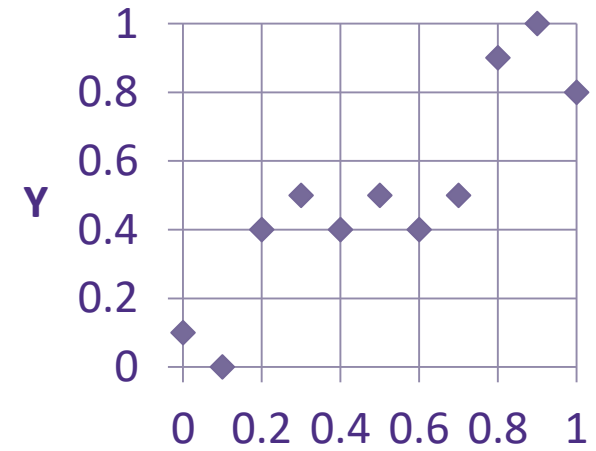
$$Y = 10 + 8 * X$$



NORMALIZATION OF A LINEAR RELATIONSHIP (4)



Normalize

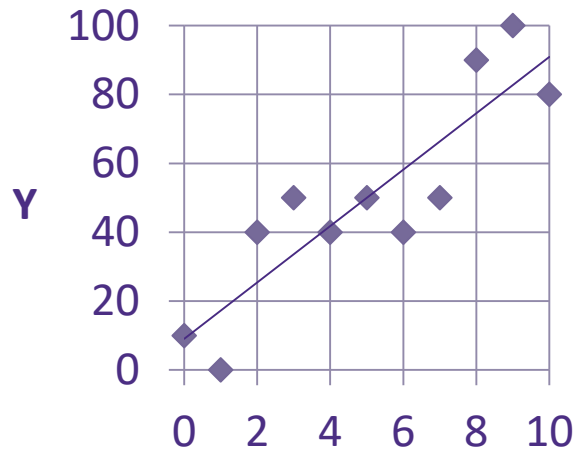


$$Y = 10 + 8 * X$$

X	Y
0	10
1	0
2	40
3	50
4	40
5	50
6	40
7	50
8	90
9	100
10	80

X	Y
0	0.1
0.1	0
0.2	0.4
0.3	0.5
0.4	0.4
0.5	0.5
0.6	0.4
0.7	0.5
0.8	0.9
0.9	1
1	0.8

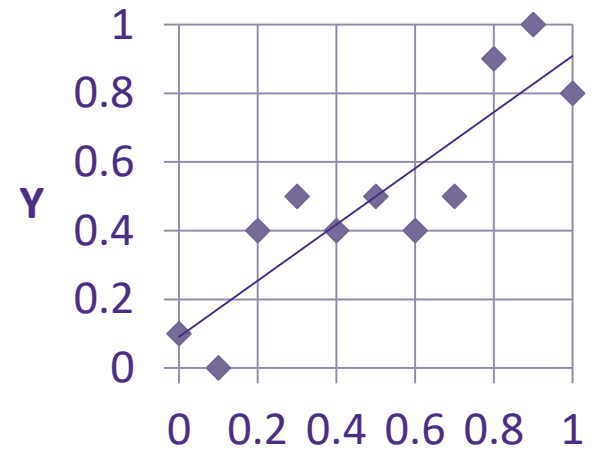
NORMALIZATION OF A LINEAR RELATIONSHIP (5)



$$Y = 10 + 8 \cdot X$$

X	Y
0	10
1	0
2	40
3	50
4	40
5	50
6	40
7	50
8	90
9	100
10	80

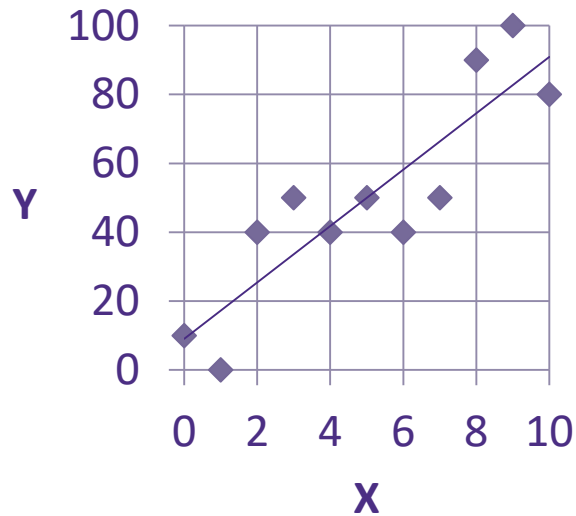
Normalize



$$Y = 0.1 + 0.8 \cdot X$$

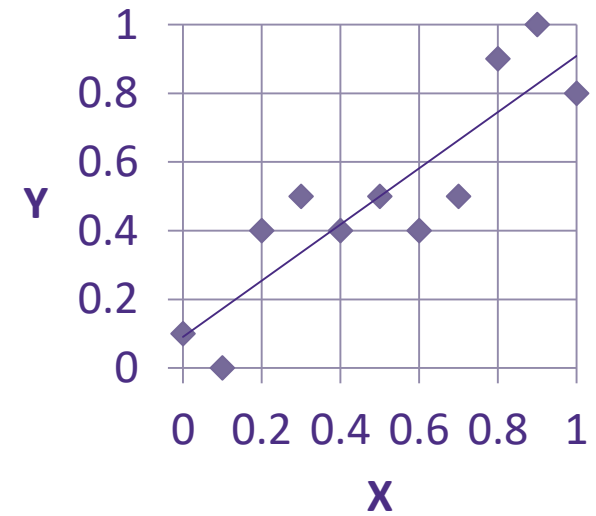
X	Y
0	0.1
0.1	0
0.2	0.4
0.3	0.5
0.4	0.4
0.5	0.5
0.6	0.4
0.7	0.5
0.8	0.9
0.9	1
1	0.8

NORMALIZATION OF A LINEAR RELATIONSHIP (6)



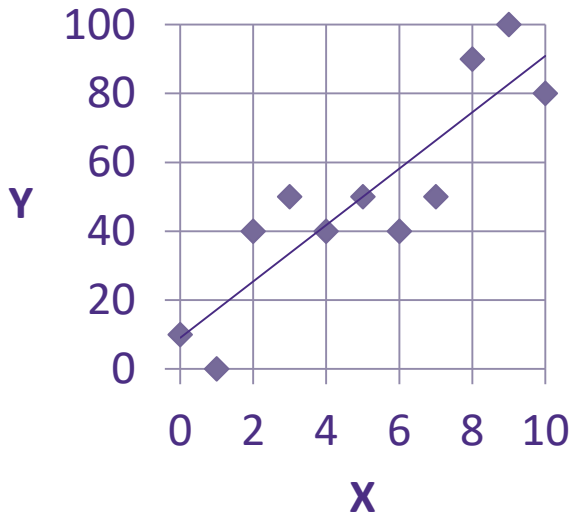
$$Y = 10 + 8 * X$$

Normalize



$$Y = 0.1 + 0.8 * X$$

NORMALIZATION OF A LINEAR RELATIONSHIP (7)



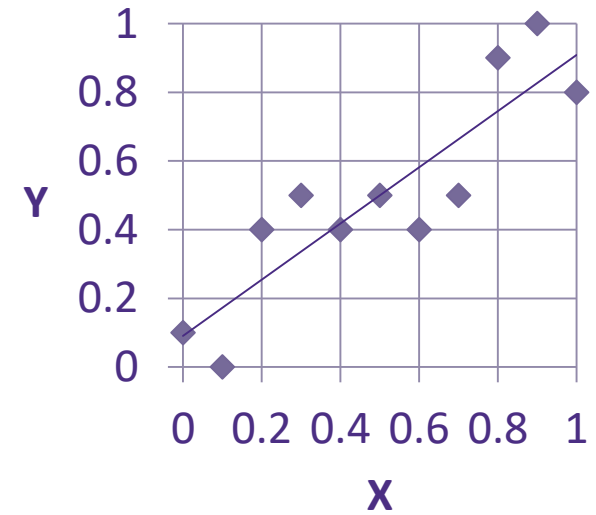
$$Y = 10 + 8 * X$$



Normalize Input
 $X = 2 \rightarrow X' = 0.2$

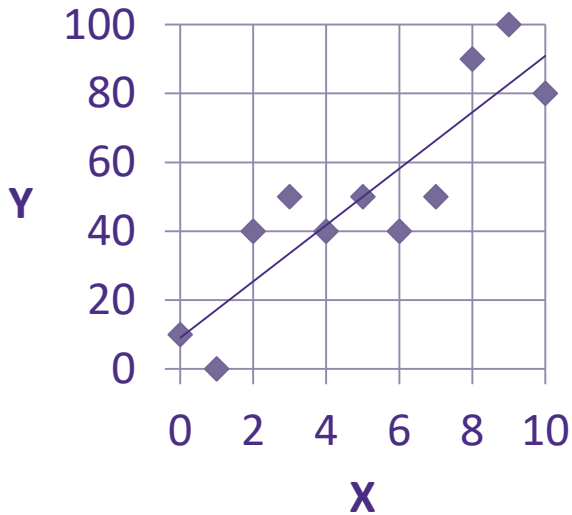
Predict Output
 $X' = 0.2 \rightarrow Y' = 0.26$

Denormalize Output
 $Y' = 0.26 \rightarrow Y = 26$



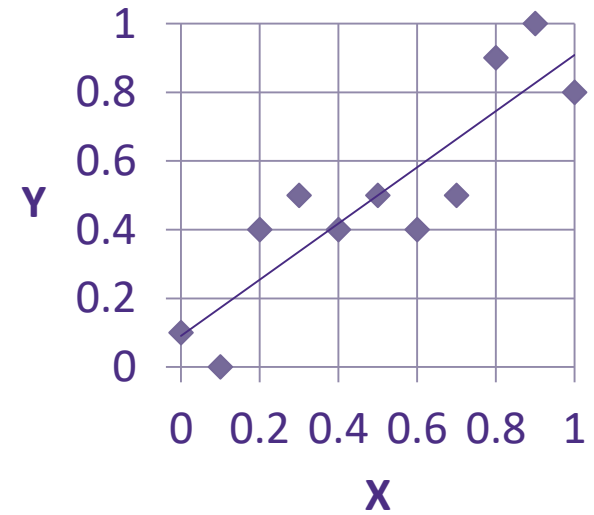
$$Y = 0.1 + 0.8 * X$$

NORMALIZATION OF A LINEAR RELATIONSHIP (8)



$$Y = 10 + 8 \cdot X$$

Normalize



$$Y' = 0.1 + 0.8 \cdot X'$$

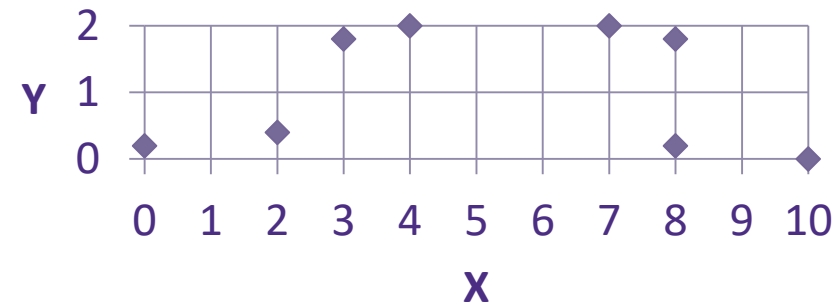
Normalize Input
 $X = 2 \rightarrow X' = 0.2$

Predict Output
 $X' = 0.2 \rightarrow Y' = 0.26$

Denormalize Output
 $Y' = 0.26 \rightarrow Y = 26$

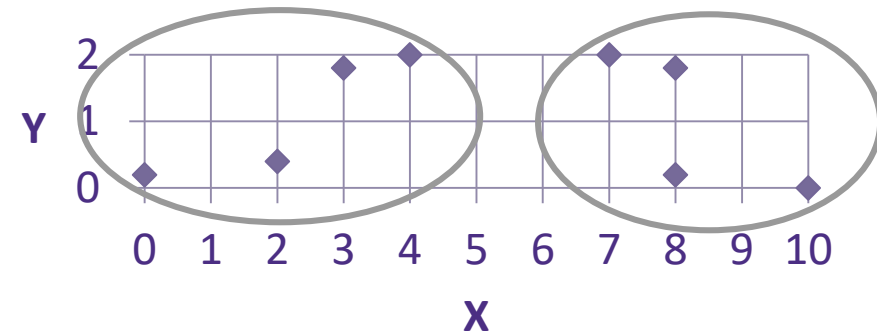
Prediction in Original Space:
 $X = 2 \rightarrow Y = 26$

NORMALIZATION OF A LINEAR RELATIONSHIP (1)



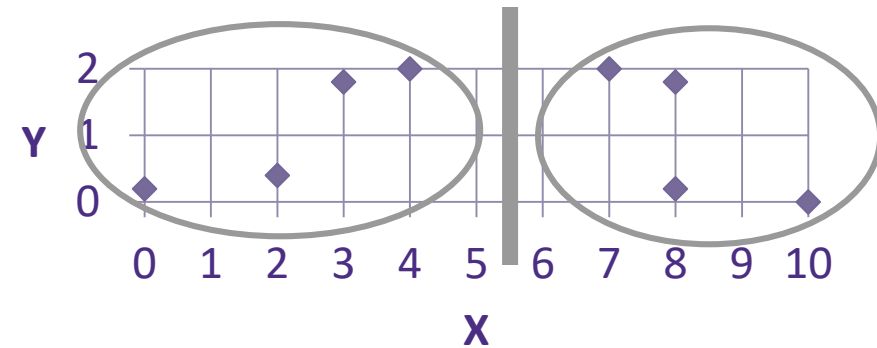
Original data in 2D:
Find 2 clusters

NORMALIZATION OF A LINEAR RELATIONSHIP (2)



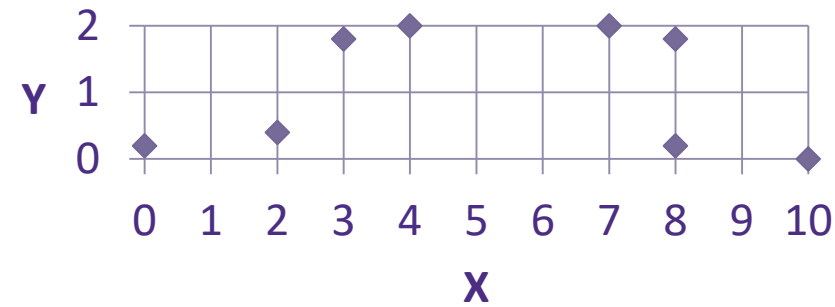
Found 2 Clusters

NORMALIZATION OF A LINEAR RELATIONSHIP (3)



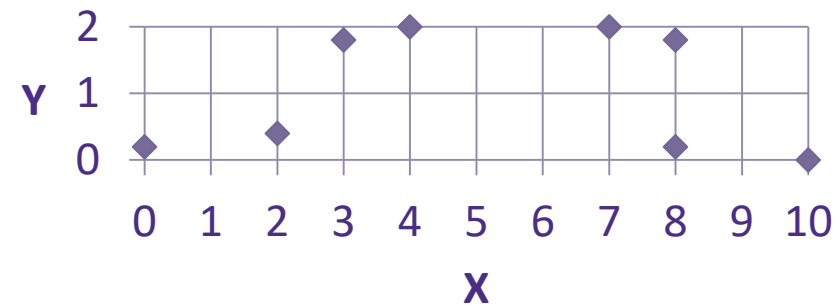
Clusters segment the image

NORMALIZATION OF A LINEAR RELATIONSHIP (4)

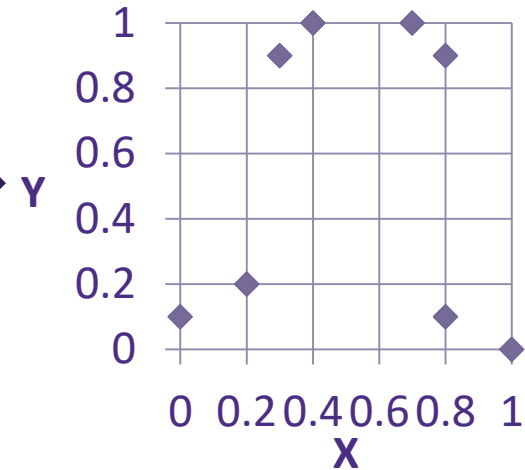


Non-normalized 2D data

NORMALIZATION OF A LINEAR RELATIONSHIP (5)

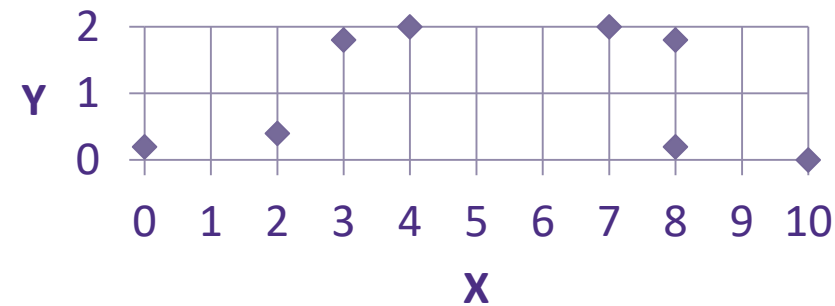


Non-normalized 2D data

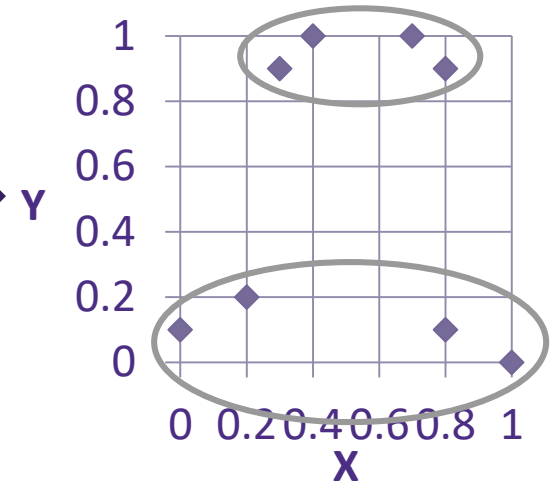


Normalize the data:
Search for 2 Clusters

NORMALIZATION OF A LINEAR RELATIONSHIP (6)

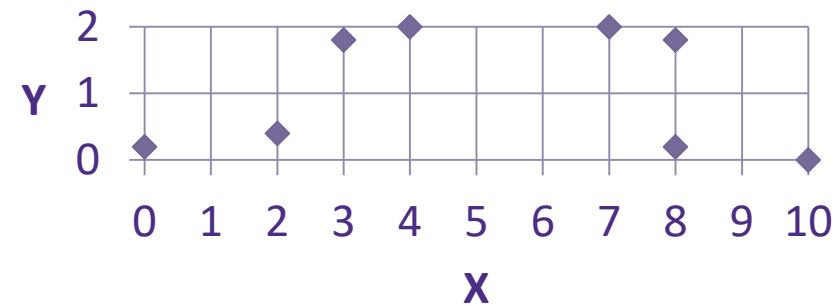


Non-normalized 2D data

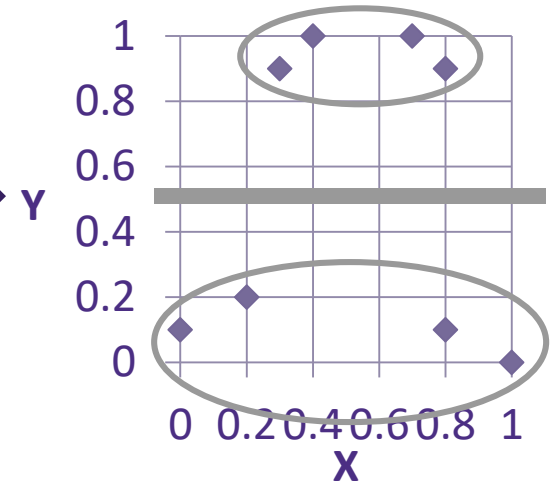


Found 2 Clusters in the normalized data

NORMALIZATION OF A LINEAR RELATIONSHIP (6)

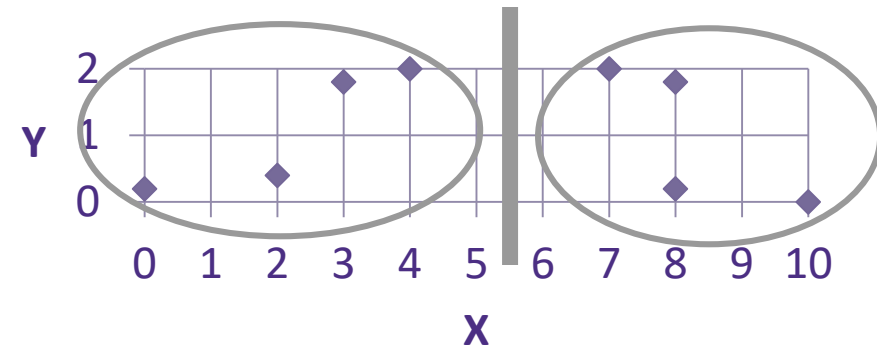


Non-normalized 2D data

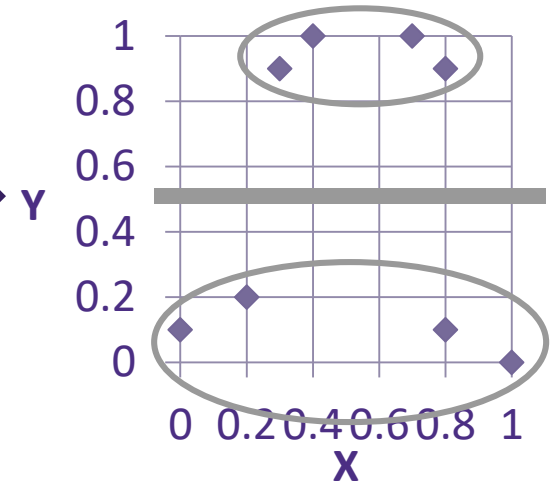


Clusters Segment the Image

NORMALIZATION OF A LINEAR RELATIONSHIP (7)



Clustering before
normalization



Clustering after
normalization

NORMALIZATION OF LINEAR AND NON-LINEAR OUTCOMES

- > Non-linear (Normalization can change outcome):
 - K-Means
 - Neural Net
- > Linear (Normalization should not change outcome):
 - Logistic Regression
 - Linear Regression
 - Mixture of Gaussians
- > <https://en.wikipedia.org/wiki/Linearity>
- > https://en.wikipedia.org/wiki/Linear_function



Normalization in Clustering



IN-CLASS EXERCISE

Normalization in K-Means

- > Download L07-2-KMeansNorm_Incomplete.py from Canvas and load into Spyder.
- > Run the script: Some results will be wrong
- > Add code to normalize each input dimension
- > Add code to de-normalize the output
- > Specifically, replace all lines that say: “Replace this line with code”.
- > Run the script: Results should be correct



IN-CLASS EXERCISE

1. L07-2-KMeansNorm_Incomplete.py
 - a. Get mean and standard deviation of point dimensions. Use the np.mean and np.std functions
 - b. Z-Normalize points and centroid guesses based on distribution of points
 - c. Let the KMeans function determine the labels and the centroids in normalized space
 - d. De-normalize the centroids
 - e. Return the labels and the de-normalized centroids



IN-CLASS EXERCISE

2. Answer the following questions

- a) What is the single most obvious difference between the distributions of the first and second dimensions?
- b) Does separation of clusters in Test 1 occur along the x, y, or both dimensions? Why?
- c) Does separation of clusters in Test 2 occur along the x, y, or both dimensions? Why?
- d) Does separation of clusters in Test 3 occur along the x, y, or both dimensions? Why?
- e) Does separation of clusters in Test 4 occur along the x, y, or both dimensions? Why?



IN-CLASS EXERCISE

3. Why is normalization important in K-means clustering?
4. How do you encode categorical data in a K-means clustering?
5. Why is clustering un-supervised learning as opposed to supervised learning?

