

دورهی آموزشی «علم داده»

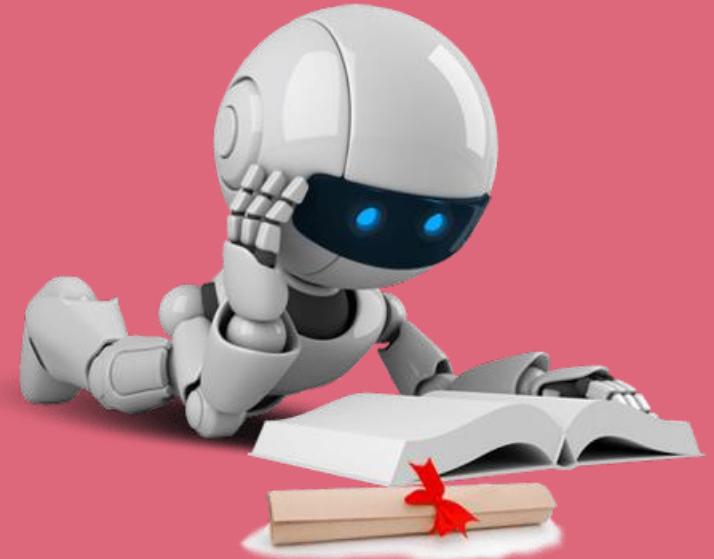
Data Science Course

جلسه بیست و سوم - (بخش اول)

نکاتی در خصوص خوشه‌بندی و استانداردسازی ویژگی‌ها



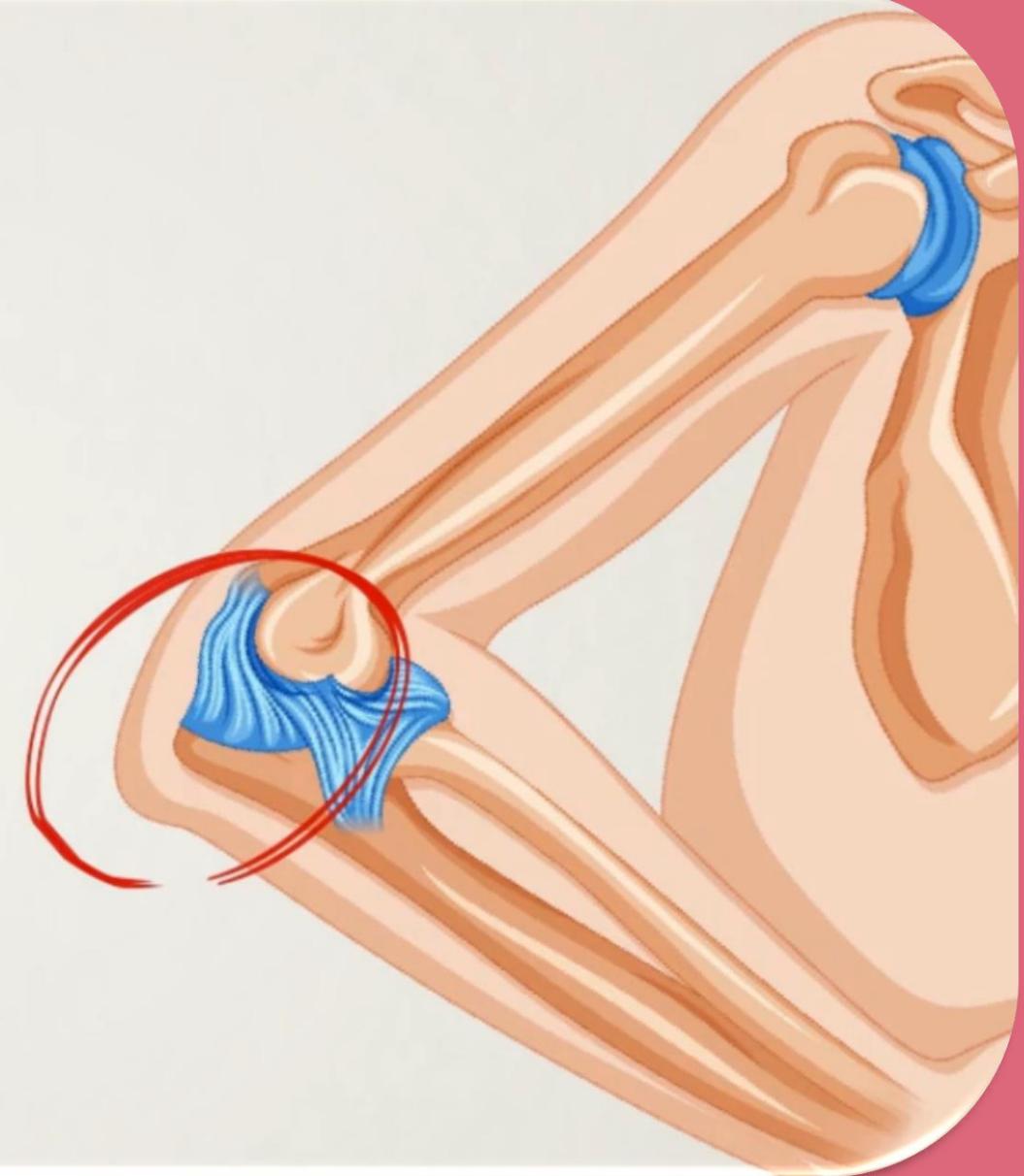
مدرس: محمد فزونی
عضو هیئت علمی دانشگاه گنبدکاووس

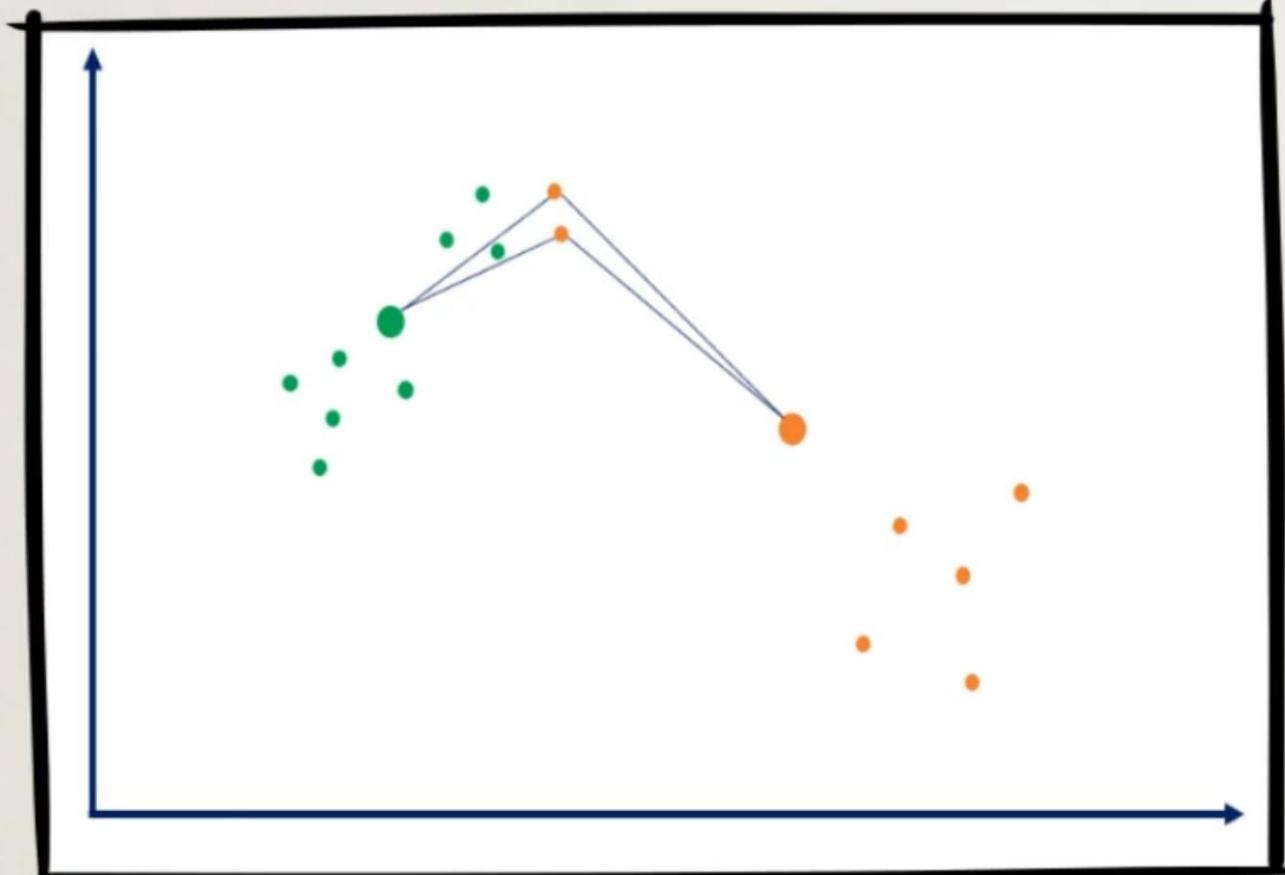


SELECTING THE NUMBER OF CLUSTERS



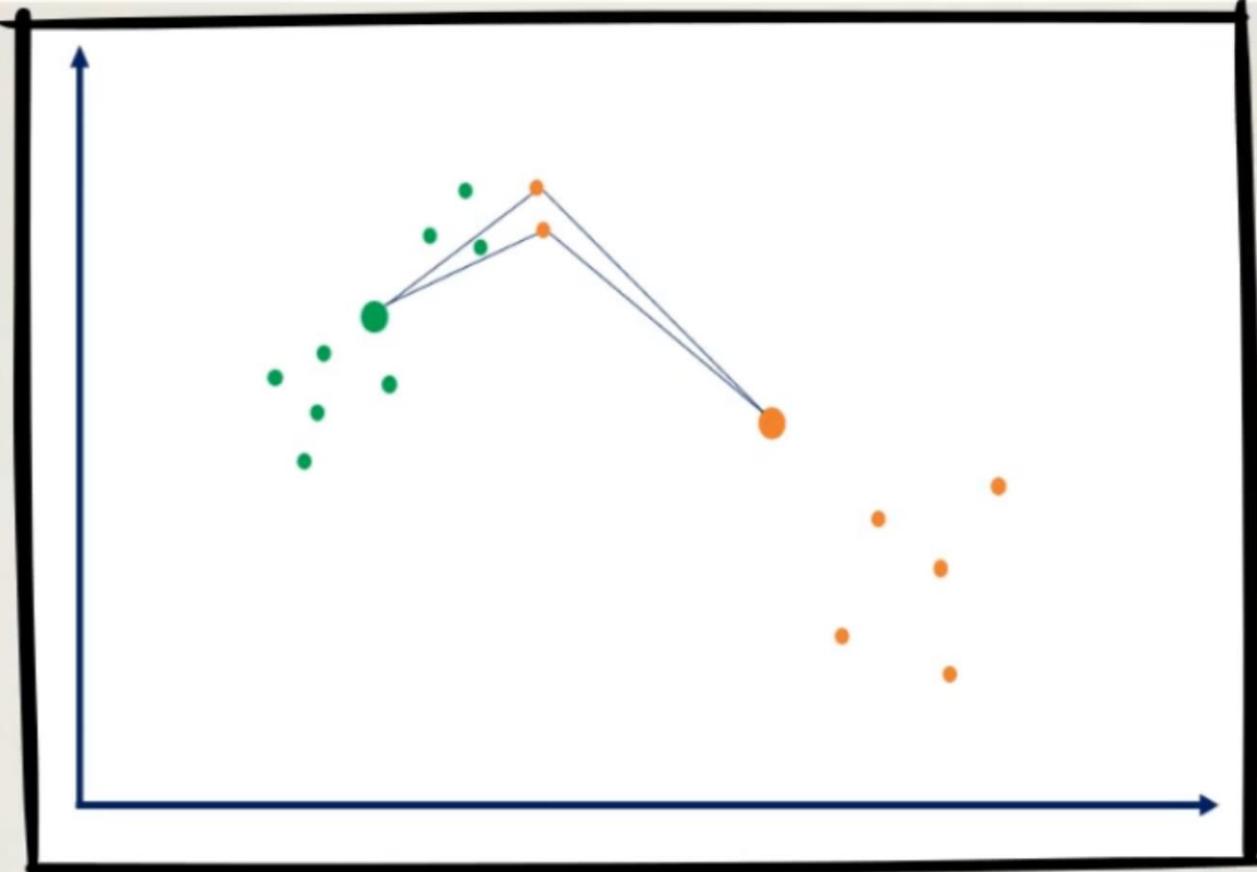
THE ELBOW METHOD



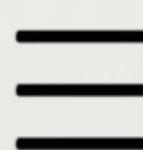


1) minimizing the distance
between points in a cluster

2) maximizing the distance
between clusters



1) minimizing the distance
between points in a cluster

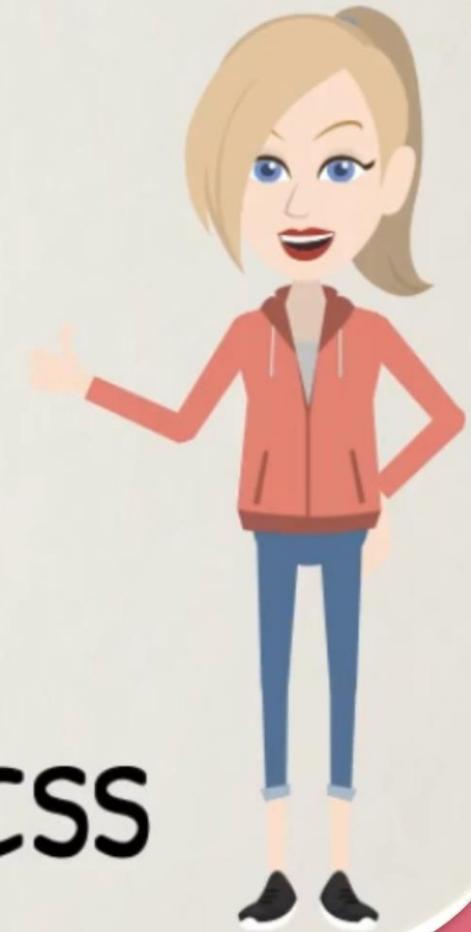


2) maximizing the distance
between clusters

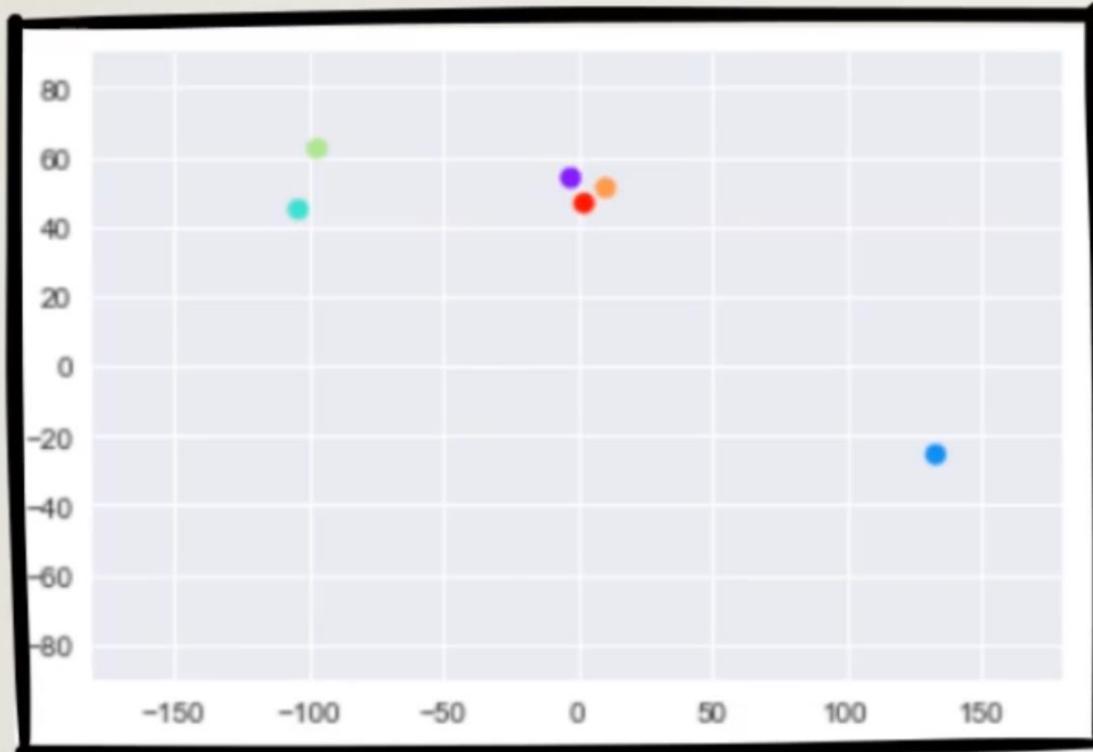
distance between points in a cluster



'within-cluster sum of squares', or WCSS



WCSS



observations: 6

clusters: 3

WCSS = 0

WCSS

observations: 1,000,000

clusters: 1,000,000

WCSS = 0 = min

observations: 1,000,000

clusters: 1

WCSS = max

MIDDLE GROUND?

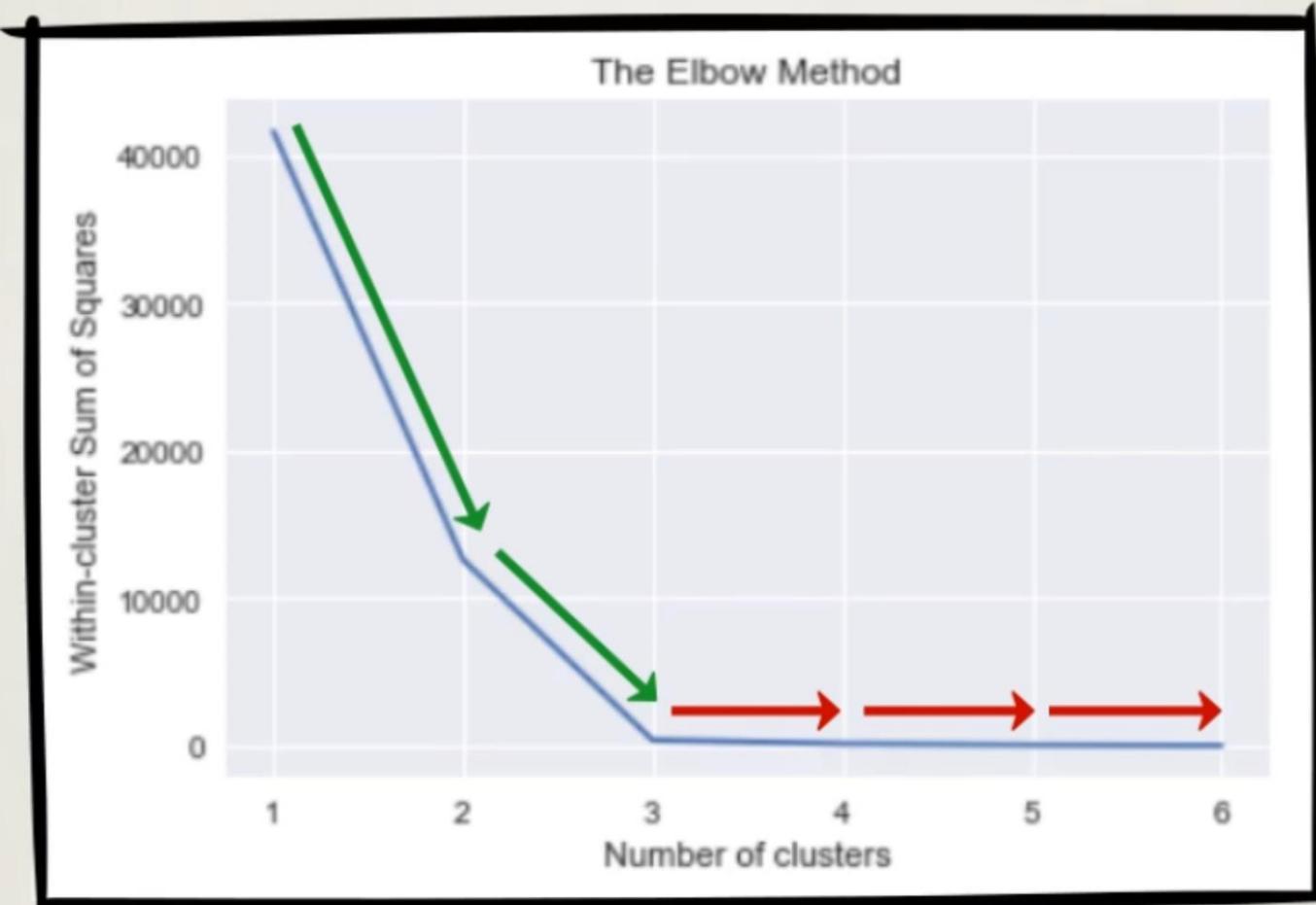
WCSS MIDDLE GROUND?

observations: **N**

clusters: **SMALL**

WCSS = LOW

THE ELBOW METHOD



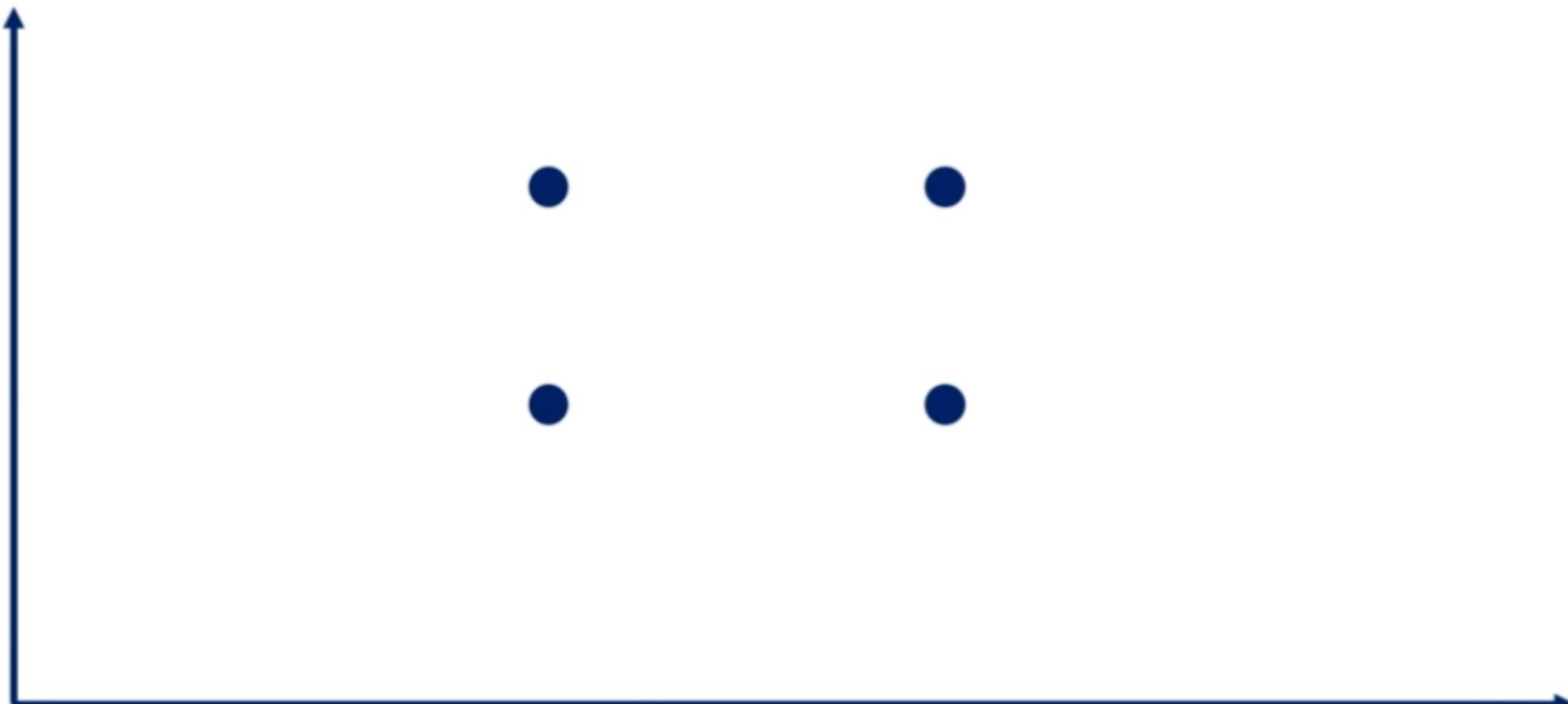
$$\text{WCSS}(k) = \sum_{j=1}^k \sum_{\mathbf{x}_i \in \text{cluster } j} \|\mathbf{x}_i - \bar{\mathbf{x}}_j\|^2,$$

where $\bar{\mathbf{x}}_j$ is the sample mean in cluster j

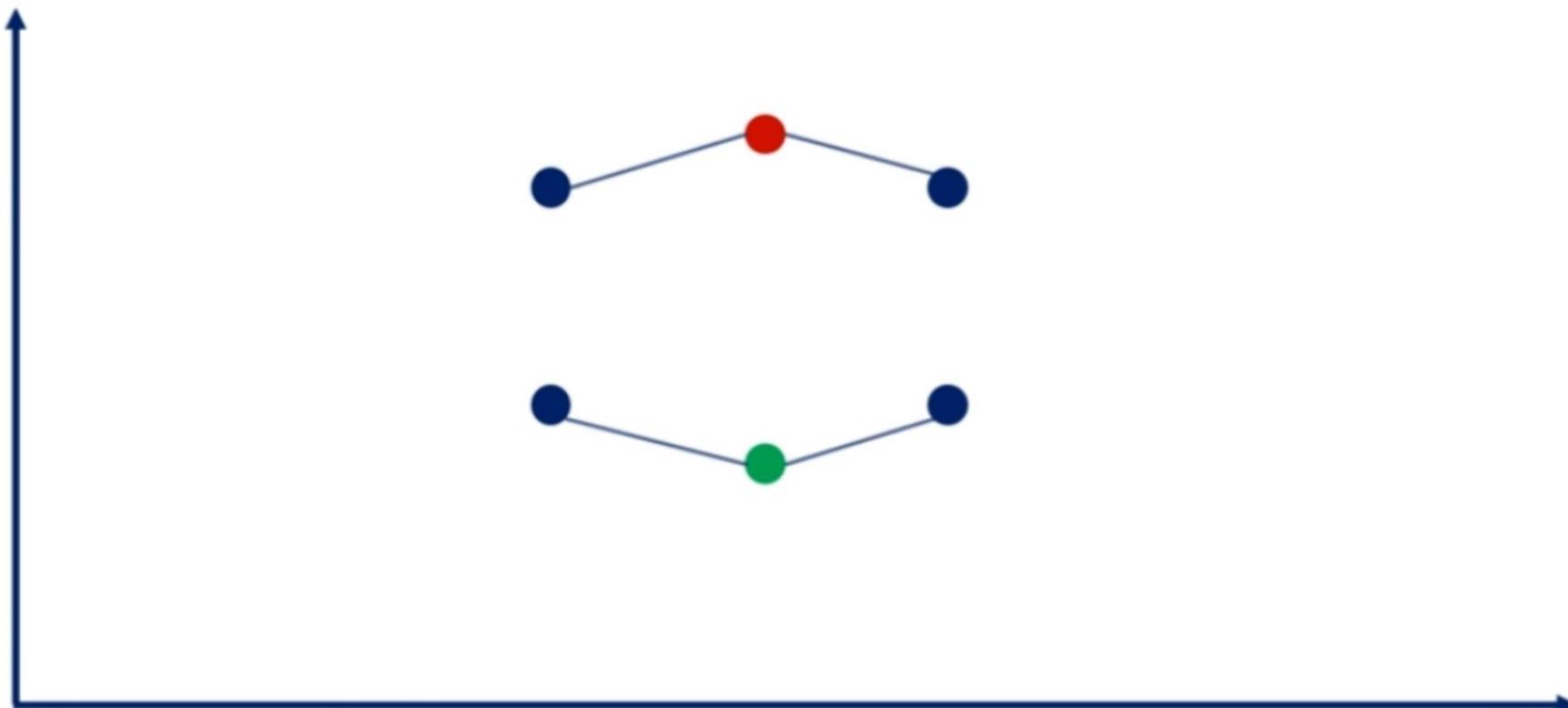
Let's re-visit

PROS AND CONS OF K-MEANS

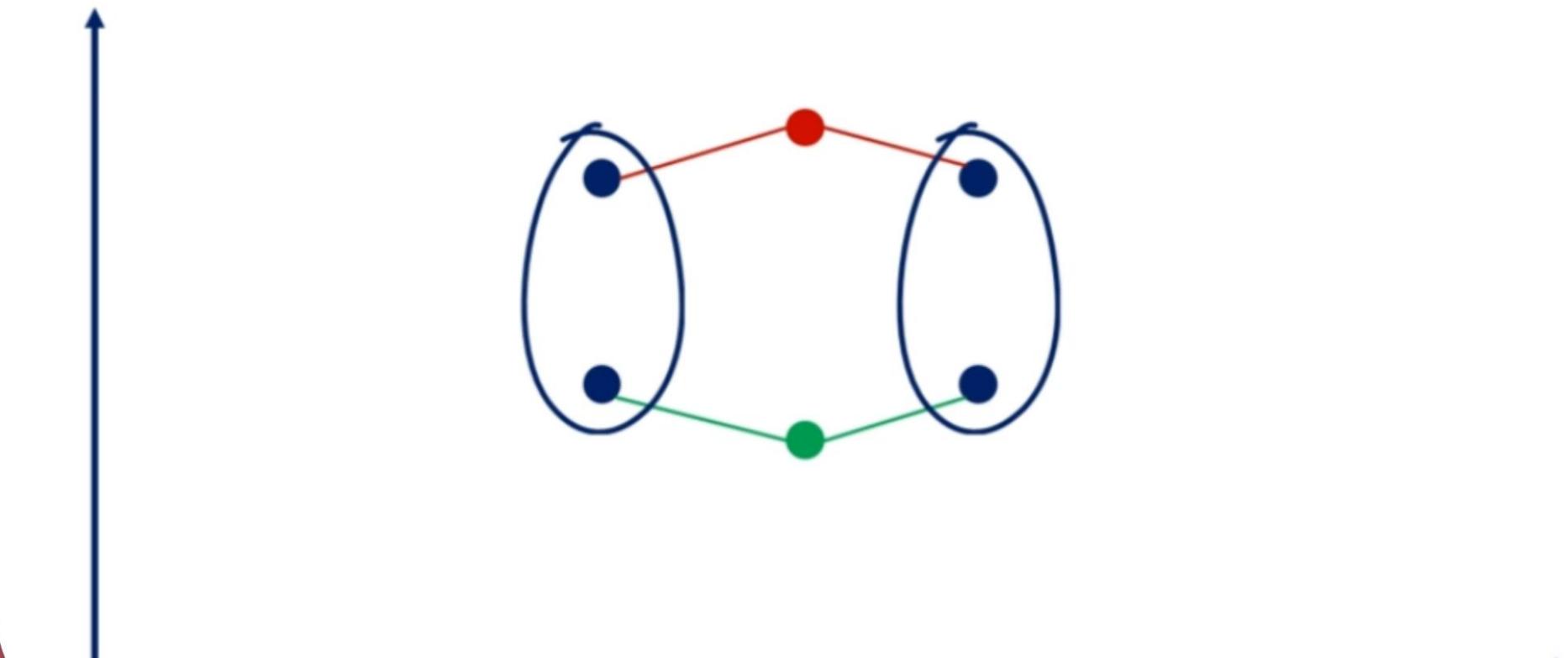
PROS AND CONS OF K-MEANS



PROS AND CONS OF K-MEANS



PROS AND CONS OF K-MEANS



What's the remedy

This algorithm ensures a smarter initialization of the centroids and improves the quality of the clustering. Apart from initialization, the rest of the algorithm is the same as the standard **K-means** algorithm. That is **K-means++** is the standard **K-means** algorithm coupled with a smarter initialization of the centroids.

Spherical

Elliptics

```
In [14]: plt.scatter(clusters_new['Satisfaction'],clusters_new['Loyalty'],c=clusters_new['Cluster'])  
plt.xlabel('Satisfaction')  
plt.ylabel('Loyalty')
```

```
Out[14]: Text(0,0.5,'Loyalty')
```

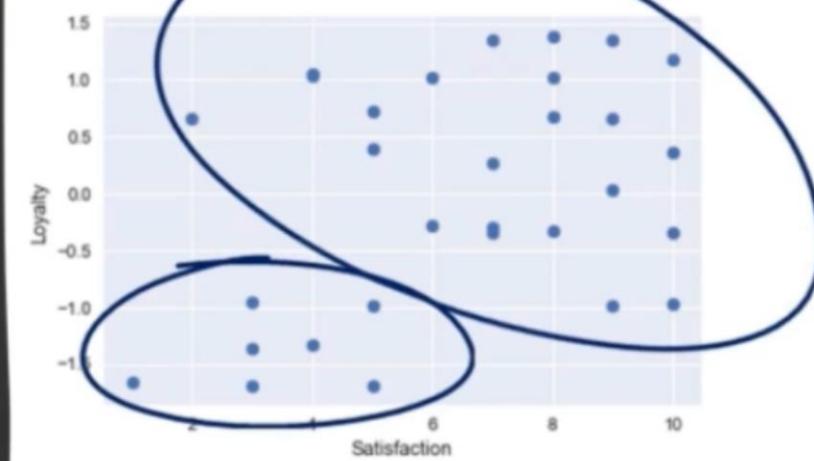


Alienated

The everything
else cluster

```
plt.scatter(data['Satisfaction'],data['Loyalty'])  
plt.xlabel('Satisfaction')  
plt.ylabel('Loyalty')
```

```
Text(0,0.5,'Loyalty')
```



CONS

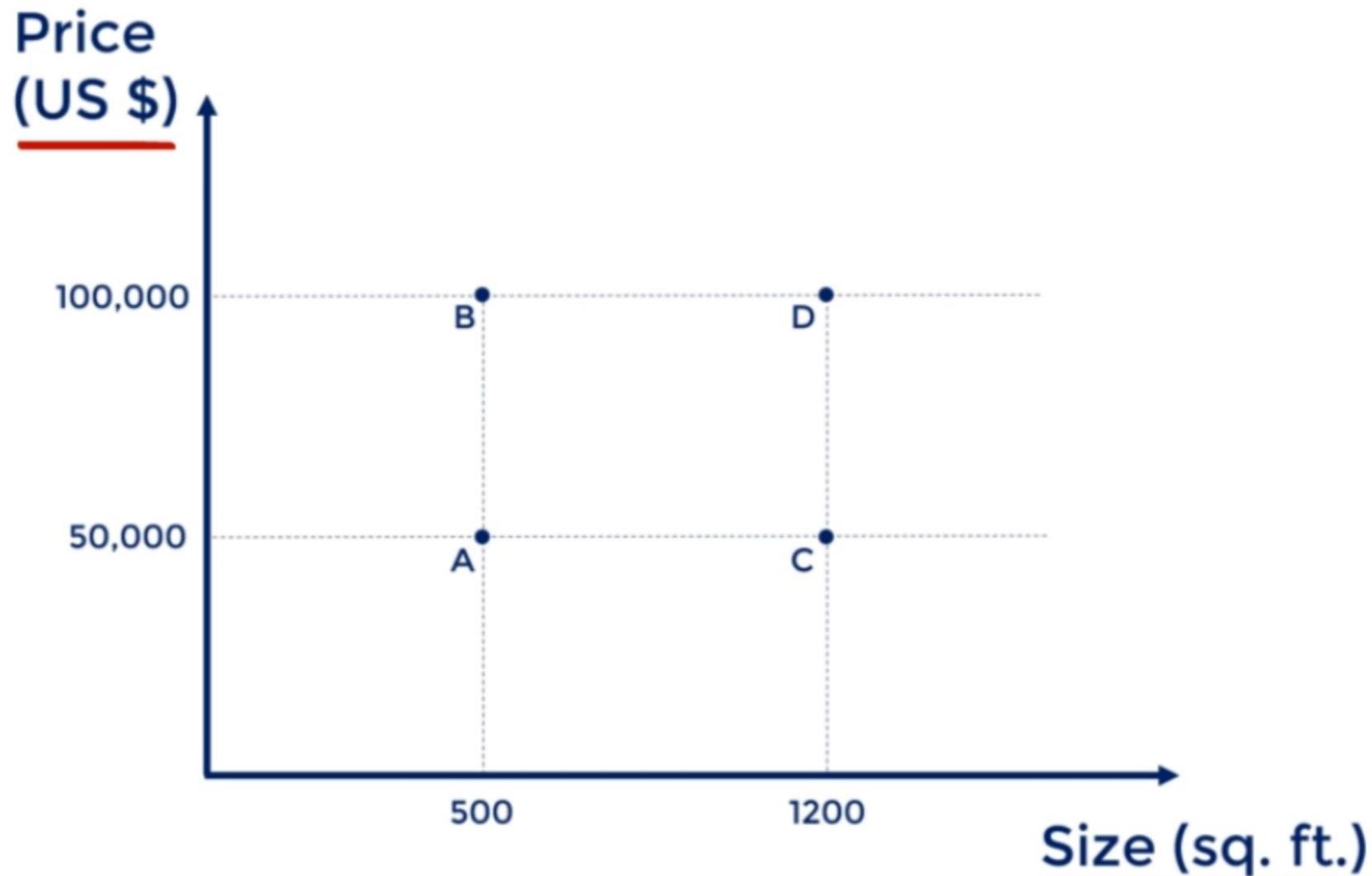
1. We need to pick K
2. Sensitive to initialization
3. Sensitive to outliers
4. Produces spherical solutions
5. Standardization

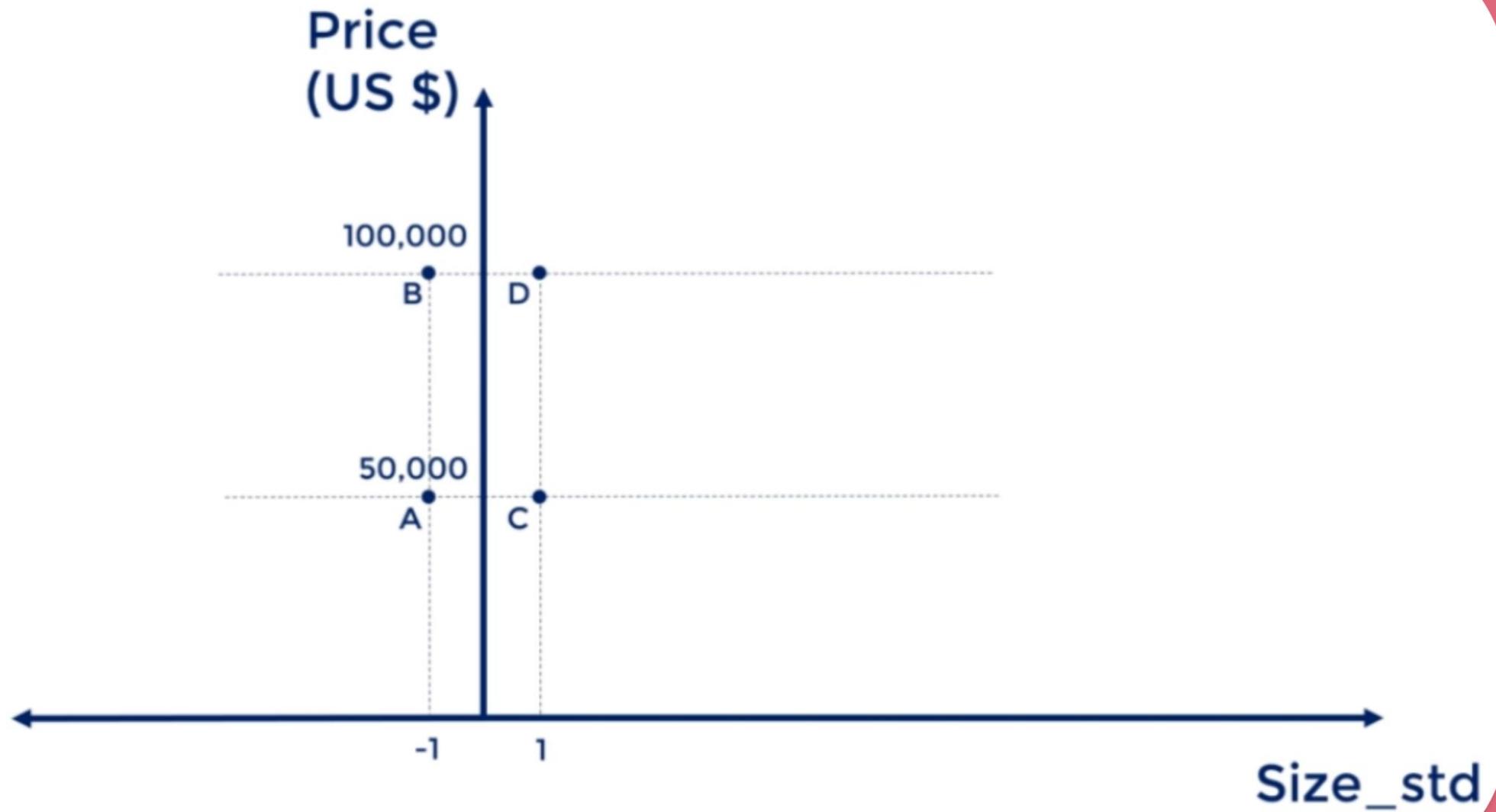
REMEDIES

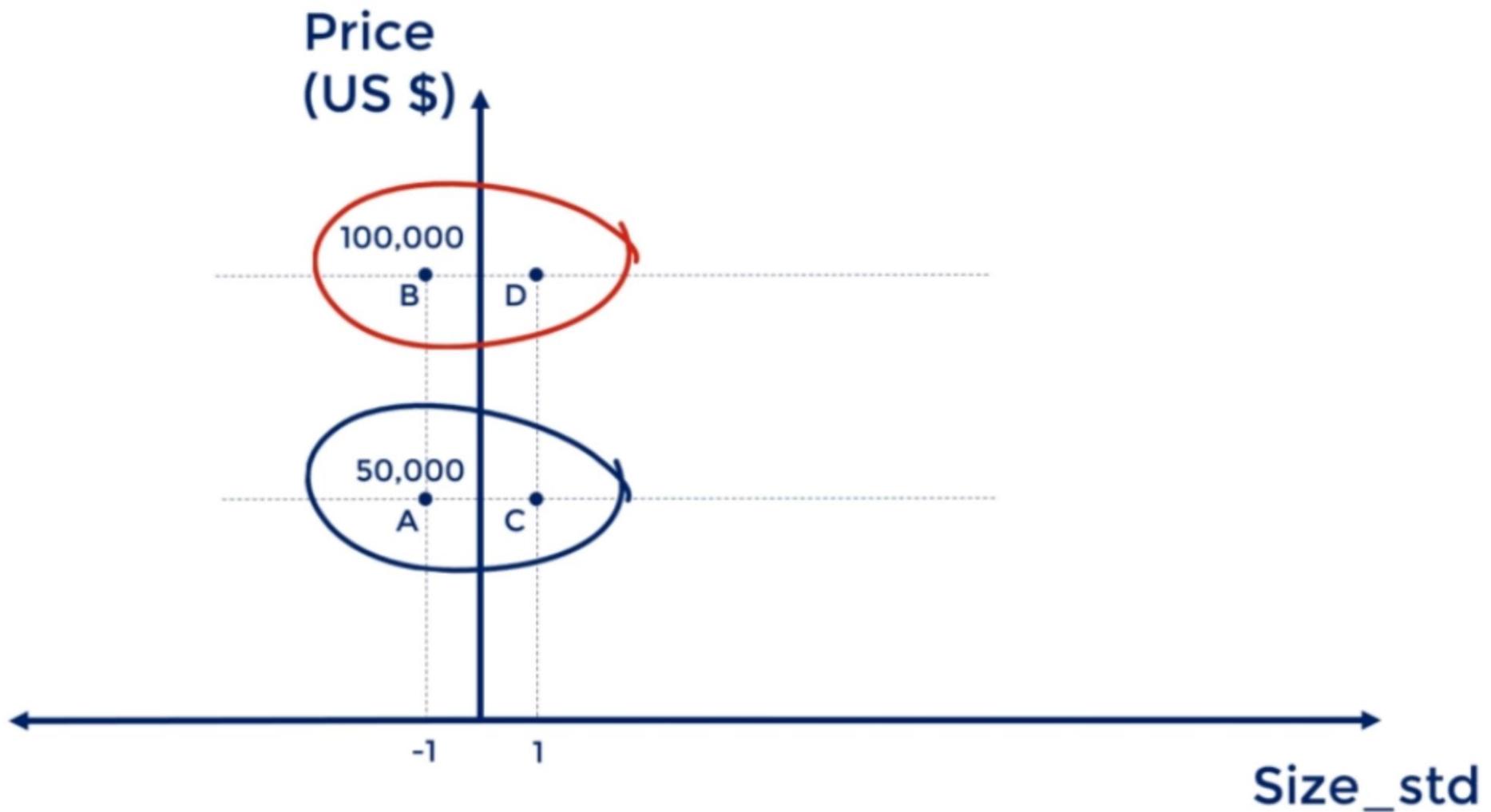
1. The Elbow method
2. k-means++
3. Remove outliers

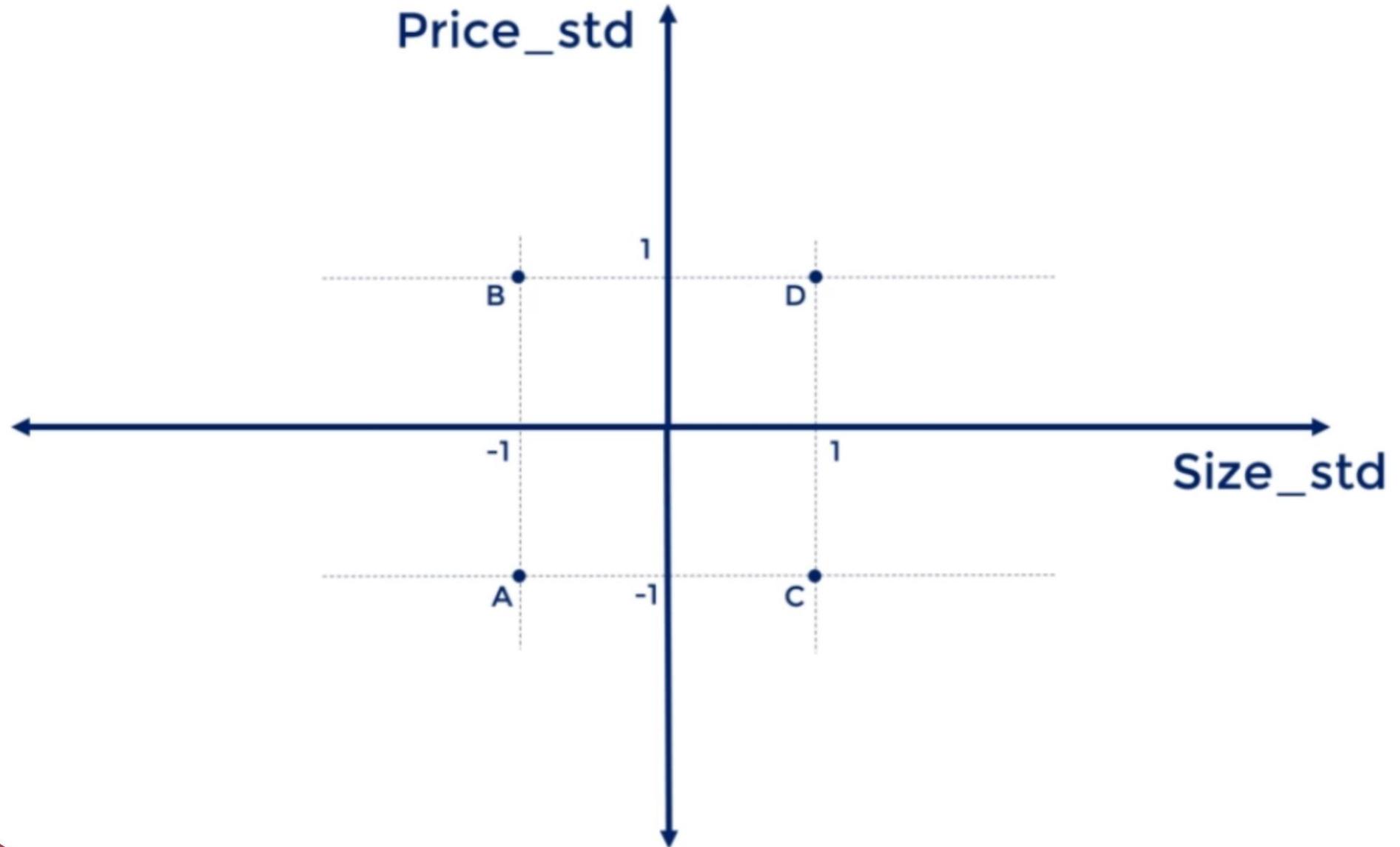
TO STANDARDIZE OR TO NOT STANDARDIZE?

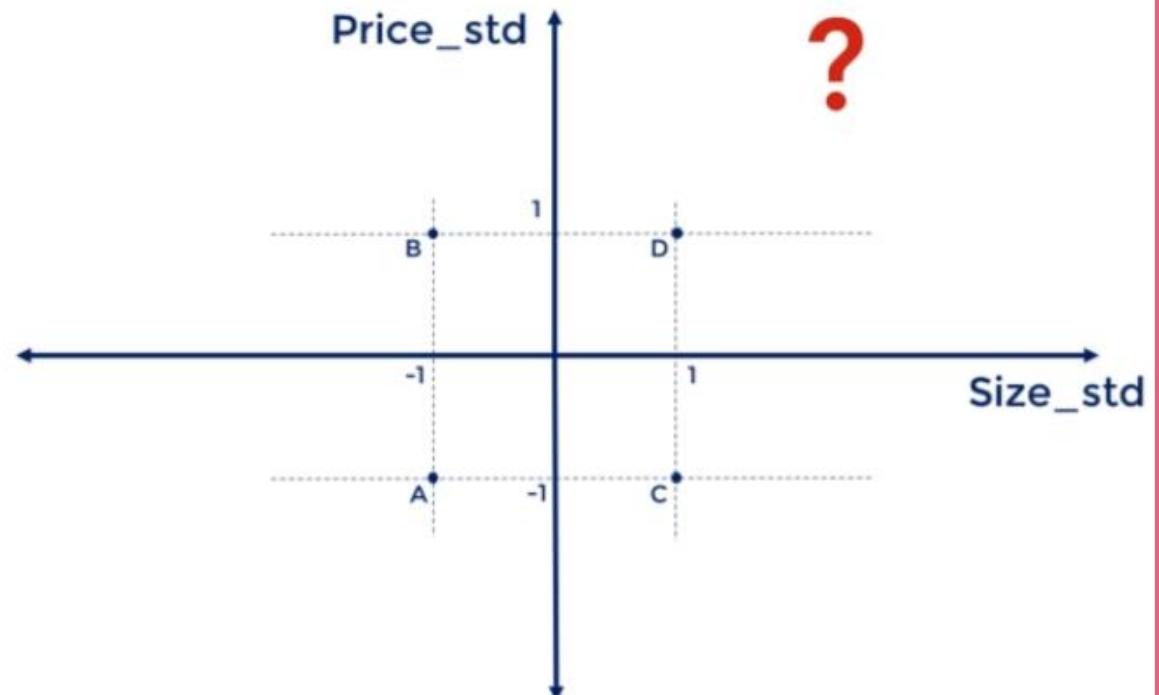
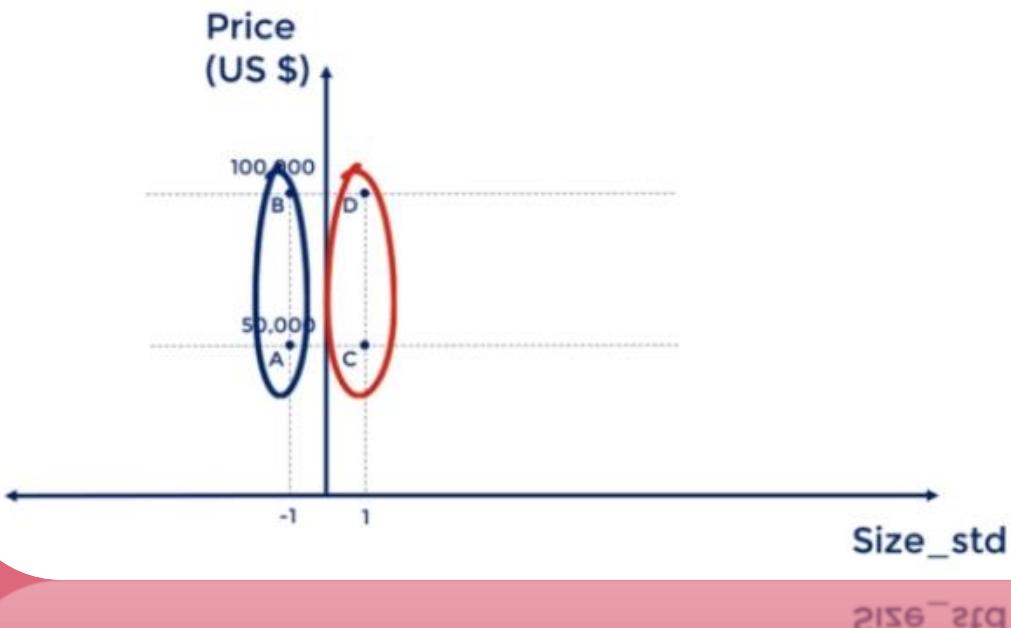
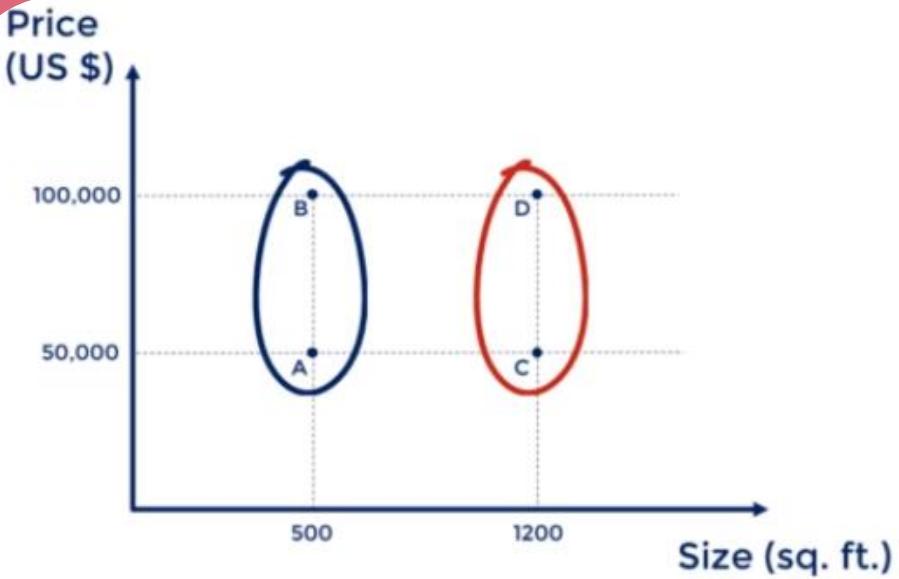


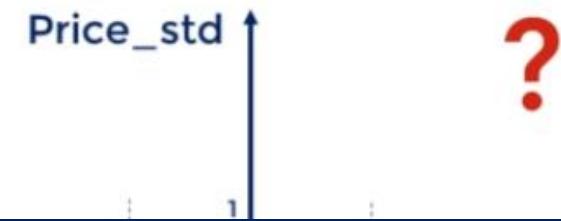
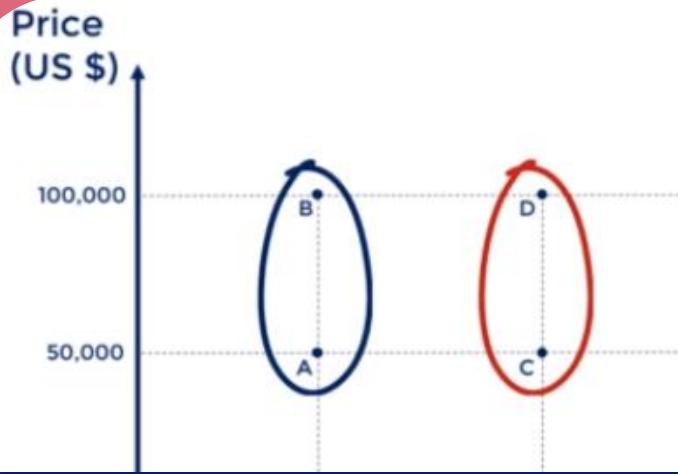




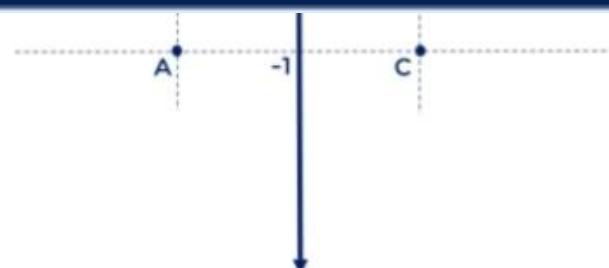
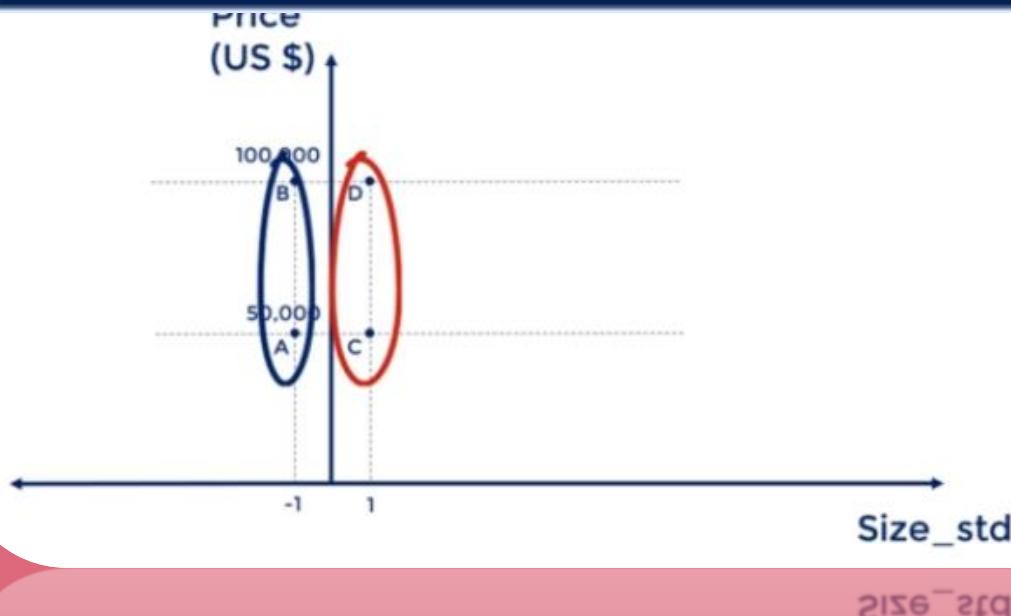


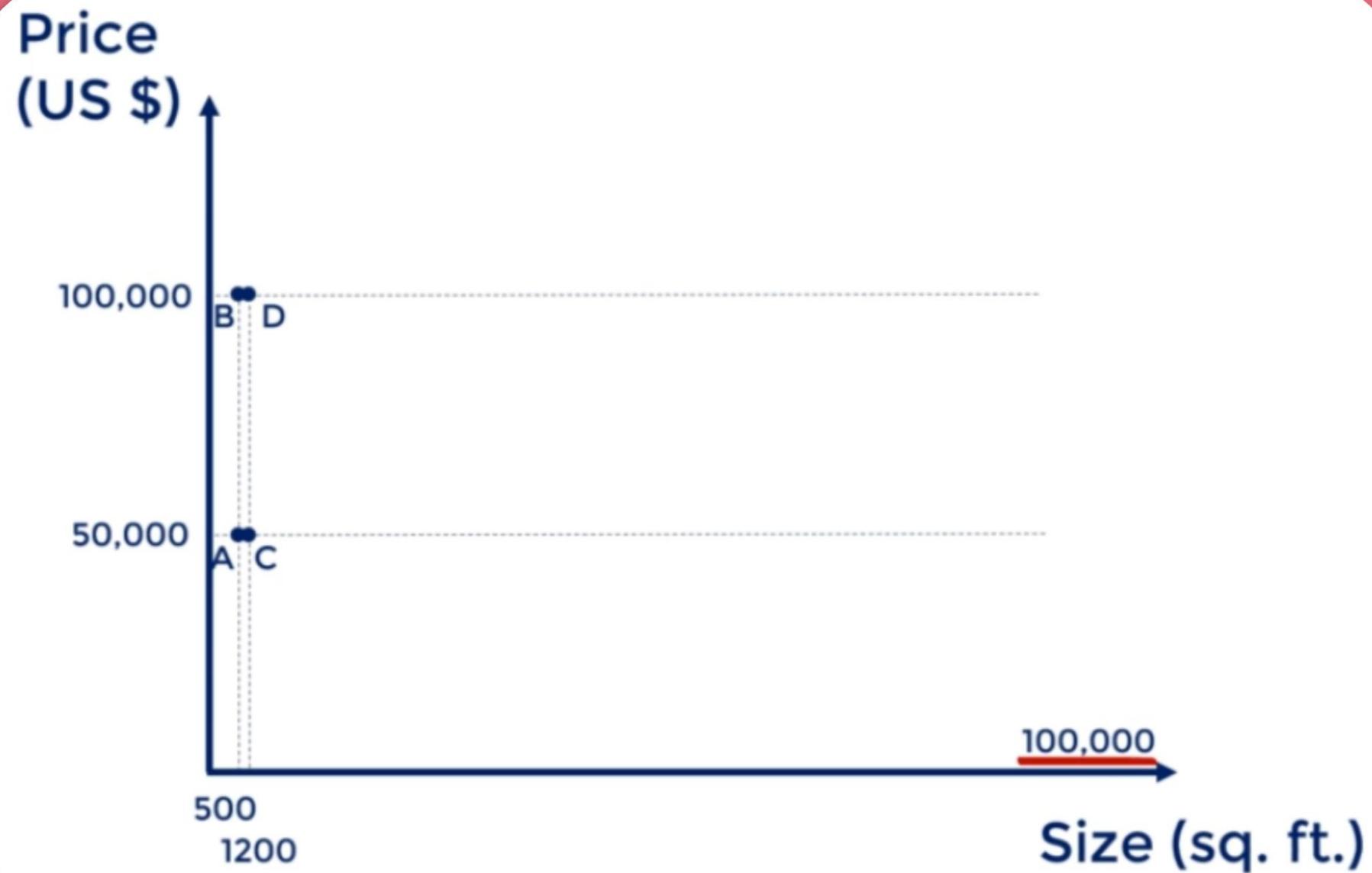


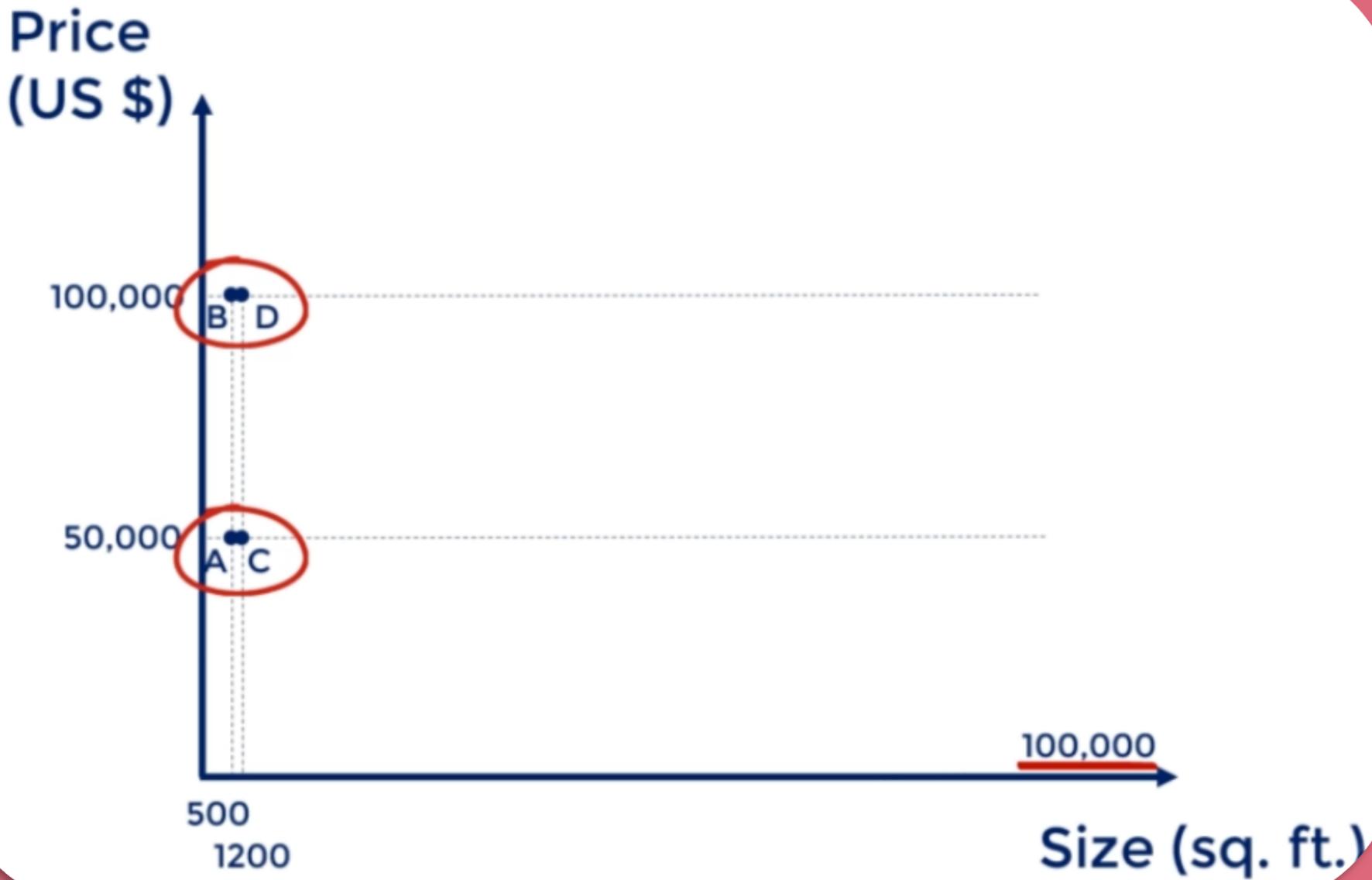




The ultimate aim of standardization is to reduce the weight of higher numbers







یک سؤال؟ کجا و کی استاندارد نکنیم؟

وقتی که دقیقاً بدونیم یکی از ویژگیها از دیگری یا بقیه مهمتر هست.
مثلاً در مثال قیمت و اندازه، قطعاً قیمت خیلی مهمتر هست. پس
استاندارد نکنیم بهتره. اینجا تجربه بهتون کمک خواهد کرد که کم کم
بفهمید ویژگی مهم کدام هست



MORE IMPORTANT

Price
(US \$)

100,000

50,000

500
1200

100,000

Size (sq. ft.)

NOT
STANDARDIZED

B D

A C





Clustering here can help us identify OVB

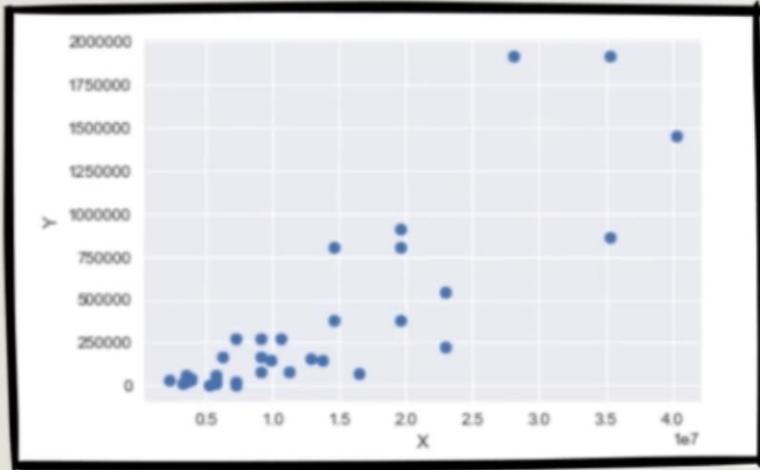
Types of analysis

- Exploratory
- Confirmatory
- Explanatory



- Get acquainted with the data
- Search for patterns
- Plan

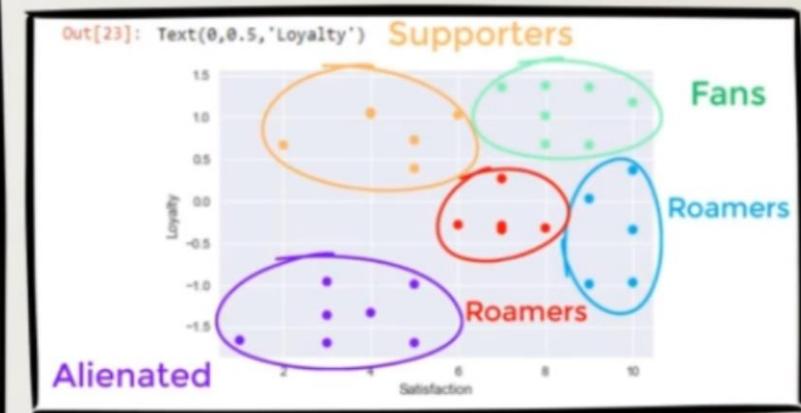
data visualization



descriptive

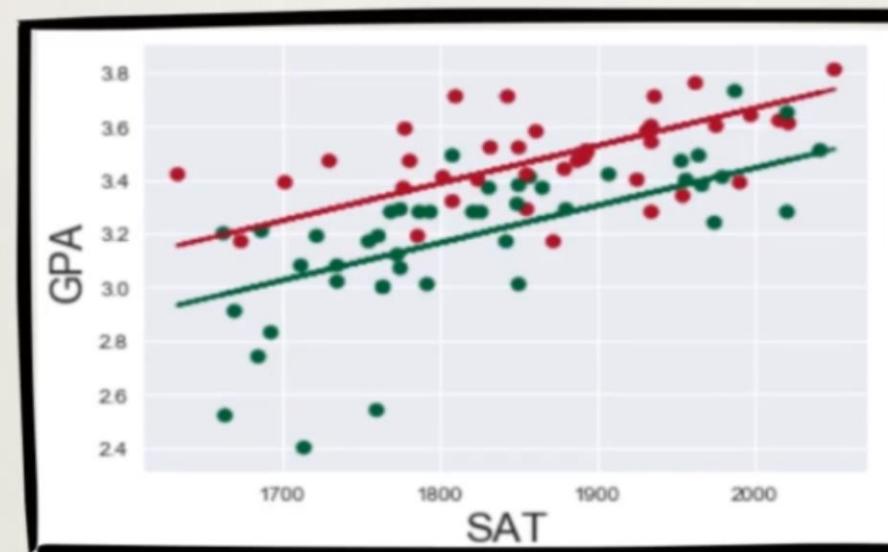
	SAT	GPA	Rand 1,2,3
count	84.000000	84.000000	84.000000
mean	1845.273810	3.330238	2.059524
std	104.530661	0.271617	0.855192
min	1634.000000	2.400000	1.000000
25%	1772.000000	3.190000	1.000000
50%	1846.000000	3.380000	2.000000
75%	1934.000000	3.502500	3.000000
max	2050.000000	3.810000	3.000000

clustering



CONFIRMATORY AND EXPLANATORY ANALYSIS

- Explain a phenomenon
- Confirm a hypothesis
- Validate previous research





Coding Time