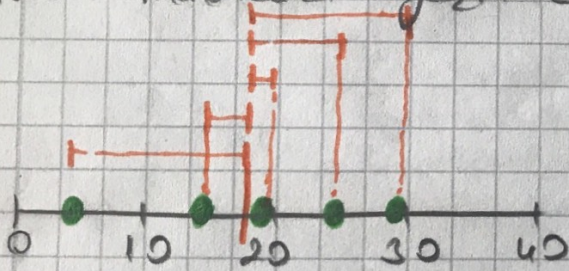


Youtube - Statquest with Josh Starmer

Covariance, Clearly Explained!!!

Varianst Hatırlatması

5 farklı moravdaki yeşil elmaların saydığını düznelim,

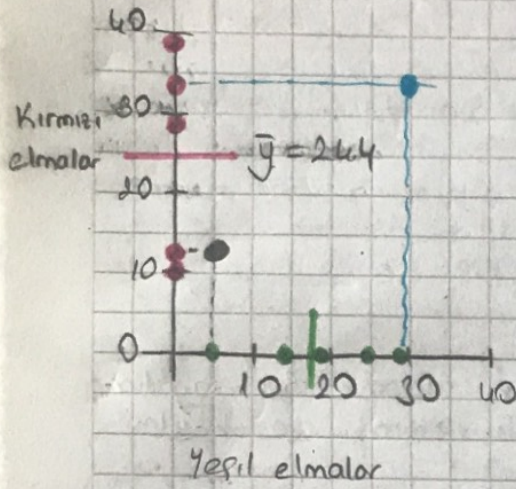


Yeşil elmaların sayısı

Estimated mean: \bar{x}

$$\text{Varianst} : \frac{\sum (x - \bar{x})^2}{n-1} = 101.8$$

✓ Yeşil elmalara ek olarak aynı manavda kırmızı elmaları da saydığımızı düşünelim.

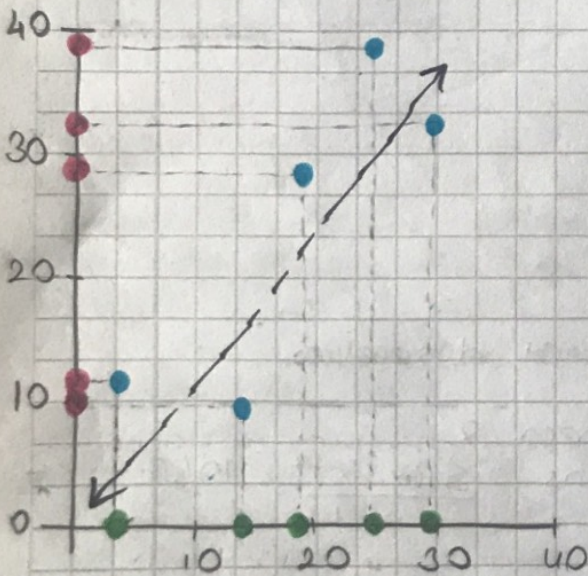


Kırmızı elmalar için varyans: $\frac{\sum (y - \bar{y})^2}{n-1} = 160.3$

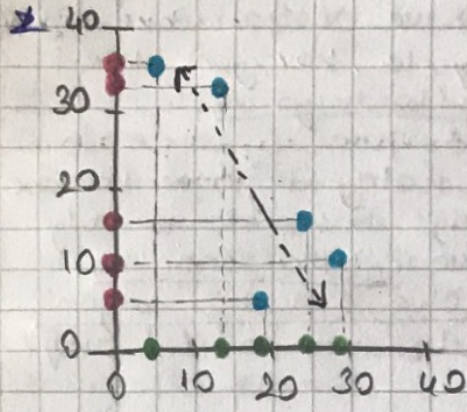
Bu ölçümler aynı manavdan alındıkları için bunları çiftler olarak (pairs) düşünebiliriz.

Örneğin; sayah ile çıktığım ölçümler aynı manavdan geliyor. Her ikisi de kendi ortalamaya değerlerinden daha az. Mayılla çıktığım da aynı manavdan geliyor çünkü her ikisi de ortalamadan daha büyük. Bu ölçümler çiftler halinde yapıldığına göre sorumuz şu: "İkili olarak alınan ölçümler bize birleşik ölçümleri söylemediği bir şey mi söyler?" **Covariance** bu soruyu cevaplamak için bir yoldur.

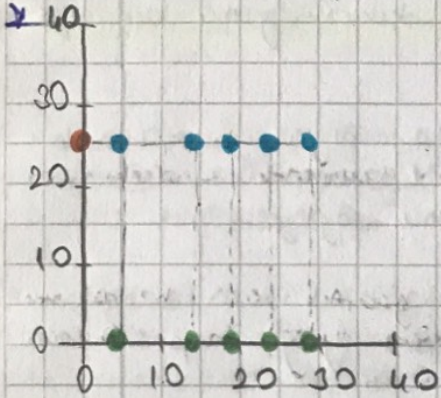
✓ Bu ölçümler aynı manavdan geldiğine göre bunları birleştirip noktalar halinde ifade edebiliriz.



Genel konuşacak olursak X ekseninde düşük değere sahip olan değerler Y ekseninde düşük değere sahip, X ekseninde büyük değere sahip olan değerler de Y ekseninde büyük değere sahip. Ortadaki line ile bu ilişkiyi ifade edebiliriz. Bu artış pozitif eğimi ifade ediyor. Yani X arttıkça Y de birlikte artıyor.



Verimiz bu şekildeyse; düşük X değerleri yüksek Y değerleriyle ilişkili, yüksek X değerleri ise düşük Y değerleriyle ilişkili. O zaman çizilecek eğri negatif bir eğime sahip olacaktır. Negatif bir trend olacaktır anlamlıdır, X'in değeri artarken Y'nin değeri azalacaktır.



Verimiz bu şekilde de olabilir. Bu durumda X'in her değeri için Y aynı değere sahip. Yani, pozitif ya da negatif X ile Y arasında herhangi bir trend yok. Y'nin bütün değerleri art olduğu durumda hangi X değeri vardı bilemeyiz bu durumda.

Benzer bir grafiği sabit bir X, değişken Y değerleri için çizdiğimizde de yazarız. 10 kırmızı elma saydıysak aynı miktarda kuru yeşil elma vardır bilemeyiz.

Kovaryansın ortasındaki ana fikir, ne tür ilişkiyi sınıflandırabiliriz. Pozitif trendleri, negatif trendleri, trendsizliği.

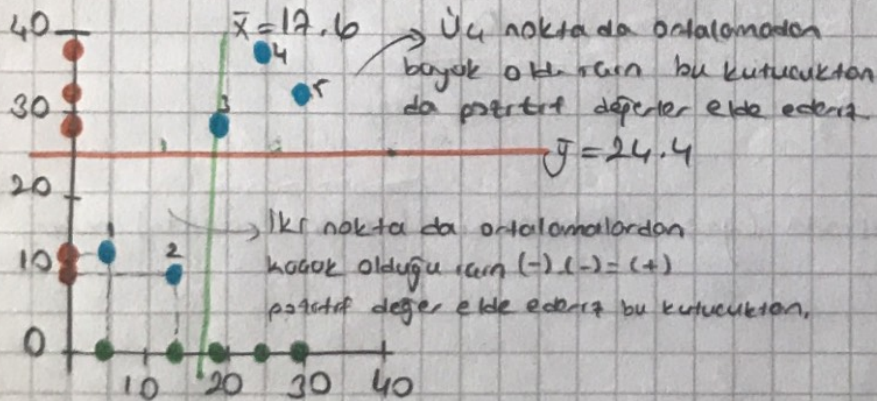
Kovaryans, kendi başına ilginç değildir. Hiçbir zaman kovaryansı hesaplamayacağız. Kovaryans, daha ilginç bir şeyin, korelasyonun bir hesaplamasıdır.

Kovaryansın Hesaplanması

$$\frac{\sum (x - \bar{x})(y - \bar{y})}{n-1}$$

1. case

Pozitif trend olan bir grafikte hesaplama yapacak olursak;



$$\begin{aligned} & \frac{(3-12.6) \times (12-24.4) + (-4.6) \times (-14.4)}{5-1} \\ & + \frac{1.6 \times 4.6}{6.4} + 55 + 155 \\ & = 116 \end{aligned}$$

11b, yani pozitif bir deęer bulduğumua rana X ve Y arasındaki ilişkinin pozitif olduğu konusunda korıyorduk. Başka bir deęerle, kovaryans deęeri pozitif olduğu rana trendi "pozitif" olarak sınıflandırıyoruz.

✗ Kovaryans deęerinin kendisini yorumlamak çok kolay deęil ve kontekte baęlı. Kovaryans deęeri bize ilişkinin dik olup olmadığı hakkında bir bilgi vermez. Sadece eğimin pozitif olduğu bilgisini verir. Daha önemlisi; kovaryans deęeri bize noktaların lineer yakın olup olmadığını da vermez. Pozitif kovaryans deęeri bize sadece ilişkinin pozitif olduğu bilgisini verir.

✗ Kovaryansı yorumlamak zor olsa da bizim rana daha ilgili bir şeyin korelasyonun bir hesaplaması basamağı.

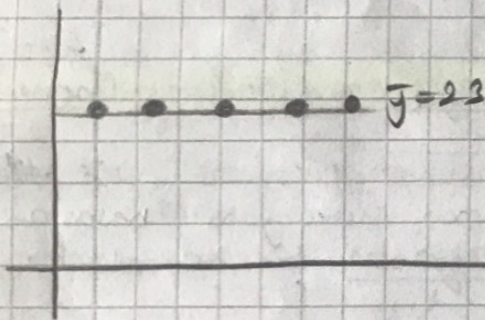
2. case

✗ X ile Y arasında negatif ilişki olan bir grafikte aynı hesaplamaları yaptığımızda sonuç -15.15 çıktı. Buradan X ve Y arasındaki ilişkinin ve çıkacak olan eğimin eğiminin negatif olduğunu söyleyebiliriz.

3. case

✗ Farklı X deęerleri rana aynı Y deęerlerini veren grafik rana hesaplaması yaparsak \bar{y} , kesim noktalarının ortasında olacağından $(y - \bar{y})$ hep 0 çıkacaktır, dolayısıyla sonuç da 0 çıkacaktır.

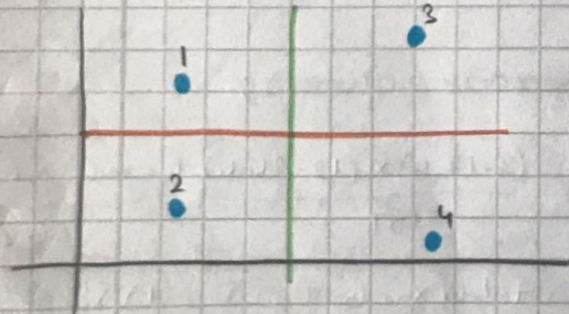
$$\frac{\sum (x - \bar{x}) (y - \bar{y})}{n-1} = \frac{(-14.6 \times 0) + 0 + 0 + 0 + 0}{4} = 0$$



4. case

✗ Verilerimiz öyle bir yerleştirebiliriz ki pozitif ve negatif deęerler birbirini götürüp 0 elde edebiliriz.

$$\frac{\sum (x - \bar{x}) (y - \bar{y})}{n-1} = \frac{-100 + 100 + 150 - 150}{4-1} = 0$$



Don iki örnekte X ve Y arasında (yeşil ve kırmızı elmalar arasında) herhangi bir ilişki olmadığını anladık.

Neden Kovaryansı yorumlamak zordur?

✓ Yeşil elmaların yine yeşil elmalarla olan kovaryansını hesaplayalım. Ortalama yine 17.6 olacaktır.

$$\frac{\sum (x - \bar{x}) \cdot (y - \bar{y})}{n-1} \text{ formülünde } y=x, \bar{y}=\bar{x} \text{ olacaktır}$$

Dolayısıyla formül;

$$\frac{\sum (x - \bar{x}) \cdot (x - \bar{x})}{n-1} = \frac{\sum (x - \bar{x})^2}{n-1} \text{ estimating varyans formülüne dönüşecektir.}$$

Bu işlemi yaptığımızda 102 sonucunu elde ettik. Bu kovaryans değeri pozitif olduğu için X'in kendisiyle olan ilişkisinin pozitif olduğunu söylüyor.

✓ Yeşil elma ölçümlerini 2 katına çıkardığımızda ve bu hesaplamaları tekrar yaptığımızda 408 bulacağız kovaryans değerini. Aradaki ilişki değişmemesine rağmen (kendisiyle olan kovaryansı hesaplıyoruz) kovaryans değeri 4 kat arttı.

Yani; kovaryans değeri verinin scale'ına, büyüklüğüne karşı duyarlı, bu da bizim onu yorumlamamızı zorlaştırıyor.

Aynı şekilde imelin noktalara yakınlığı ve uzaklığı ile kovaryans değeri arasında bir ilişki bulunmuyor.

Scale'e duyarlı olmayan bir şey hesaplamak isterseniz **Korelasyon** imdadımıza yetişiyor.

✓ Kovaryans değeri korelasyonun bir hesaplama adımı olduğu gibi **Principal Component Analysis (PCA)**'nin de bir hesaplama adımı.