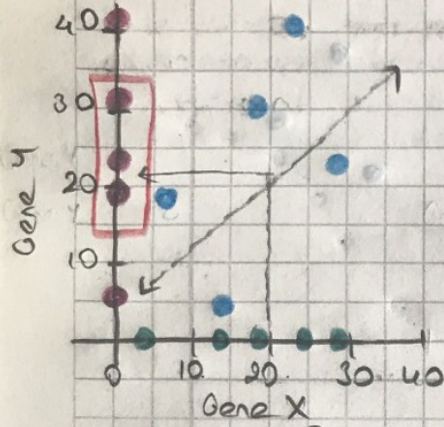


Youtube-Statquest with Josh Starmer Pearson's Correlation, Clearly Explained!!!

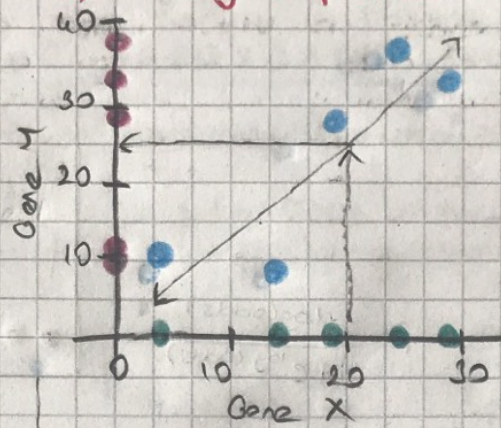


③

Eğer veri trend line'in
dan daha uzağıysa
bu durumda Gene Y daha
geniş aralığa dğer.

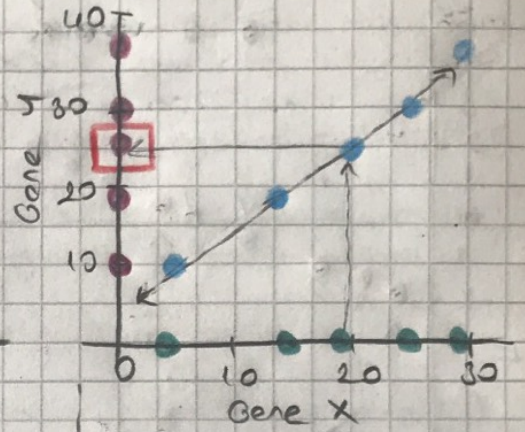
✖

Bu durumda, Gene X'in
değerleri bize Gene Y'nin
değerleri hakkında daha
az bilgi verir.
Alternatif olarak, Gene X
ve Gene Y arasındaki
ilişkinin zayıf olduğunu
söyleyebiliriz.



①

Gene X = 20 olduğunda, bu
line ile Gene Y'nin 27
aralığında olacağını tahmin
edebiliriz. Tahminlerimiz, veri
de gösterdiğim gibi trende göre
yapıyoruz.



②

Eğer veri trend line'in
daha yakınına verileri gene
X değerine göre Y,
daha küçük aralığa dğer.
Bu durumda, veriler line'a
ne kadar yakınsa X gen
bize Y hakkında daha
fazla bilgi verebilir. Alternat
tif olarak, Gene X ve Gene Y
arasındaki ilişkinin nispeten
güçlü olduğunu söyleyebiliriz.

✖

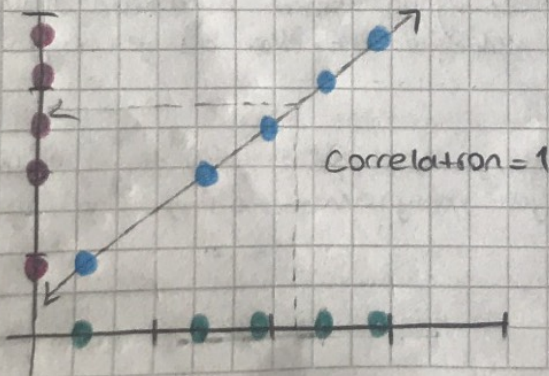
Veriler arasında zayıf
bir ilişki varsa
=
Küçük korelasyon değeri

Veriler arasında orta bir ilişki
varsa
=
Orta korelasyon değeri

Veriler arasında güçlü
bir ilişki varsa
=
Büyük korelasyon değeri

Korelasyonun Değeri

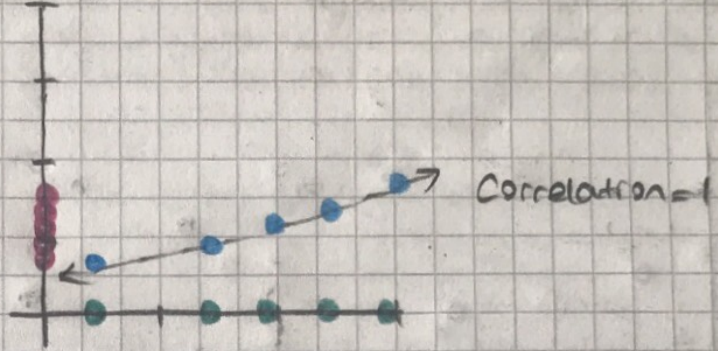
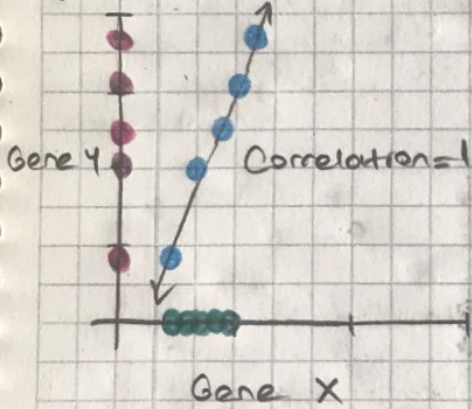
✖ Pozitif eğimli, doğu bir doğru, her ver noktasının, merkeziye göre
korelasyon = 1'dir.



Bu da bize Gene X değeri verildiğinde
Gene Y değerini çok dar bir aralıkta
tahmin edebileceğimizi anlamına gelir.

(scale)
Note: Korelasyon verinin ölçeğine bağlı değildir.

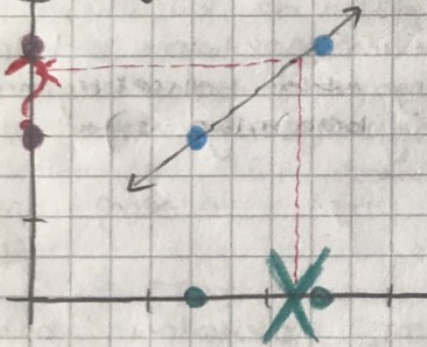
Diğer bir deyişle, verinin ölçeği ne olursa olsun, pozitif eğimli bir doğru tüm verilerden geçebildiğinde korelasyon = 1 olur. Bu, eğim büyük olsa da küçük olsa da korelasyonun 1'e eşit olabileceği anlamına geliyor.



Korelasyonun 1'e eşit olmasının ne kadar data olduğuna bir işi yok.

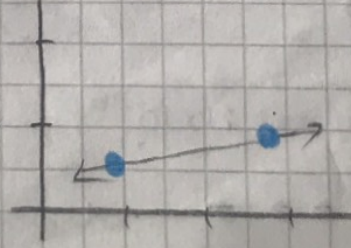
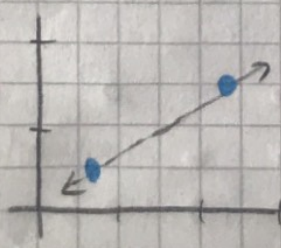
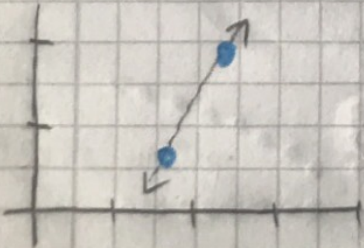
Korelasyon değeri ne zaman güvenilir?

Örneğin yalnızca 2 veri noktası varsa, iki noktayı birleştirerek pozitif eğimli doğu bir line çizebiliriz ve korelasyon = 1 olur ve bu da ilişkinin güçlü olduğunu sağlar. Ancak bu line ile yapılan tahminlere güvenmemeliyiz, çünkü elimizde çok az data var.

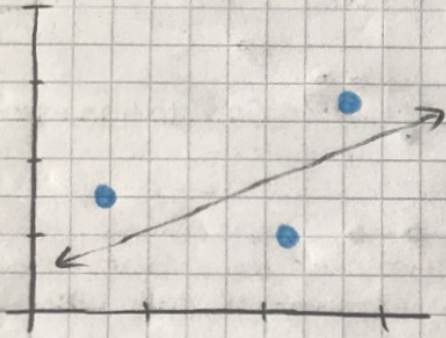


Neden küçük veri kümeleriyle yapılan korelasyonlara neden düşük güven duymalıyız?

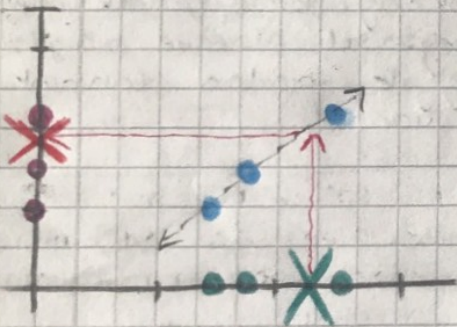
2 data point'ten oluştuğunu düşünelim. Bu noktalar nerede olursa olsun 2 random nokta arasında her açıda noktaların merkezinden geçecek şekilde de doğu çizebiliriz.



3 data pointimiz olduğunu düşünelim. 2 nokta için yaptığımız gibi 3 noktadan da da bir eğri çizebildiğimiz için korelasyon 1'dir, ancak eğer bu eğriyle yaptığımız tahminlere daha fazla güvenebiliriz. Bunun nedeni, boş bir grafikte başlayıp üzerine rastgele 3 nokta çizsek, herhangi 2 noktayı birleştirmek için da bir eğri çizmek kolay olsa da 3 noktanın hepsinden da bir eğri çizme şansımızın düşük olmasıdır.



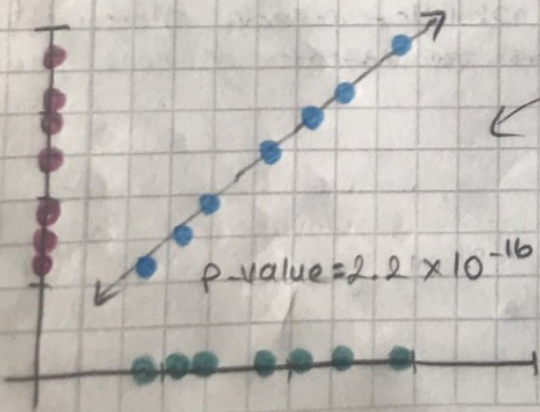
Dolayısıyla, rastgele atılmış 3 noktayı da bir eğriyle birleştirme olasılığımız çok küçüktür ve bu nedenle, gözlenen korelasyonun tesadüf sonucu göstermediğine daha fazla güvenebiliriz.



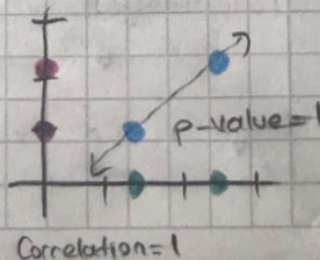
Genel olarak, ne kadar çok veriye sahipsek, doğru ile yaptığımız tahminlere o kadar güveniriz çünkü aynı sayıda rastgele yerleştirilmiş noktadan da bir eğri çizme olasılığımız her ek nokta ile daha da küçülür. (Dalgali bir eğriyle tüm noktaları birleştirebiliriz ama korelasyondan bahsettiğimiz zaman da eğriden bahsediyoruz.)

Verinin büyüklüğüyle p-value'nun ilişkisi:

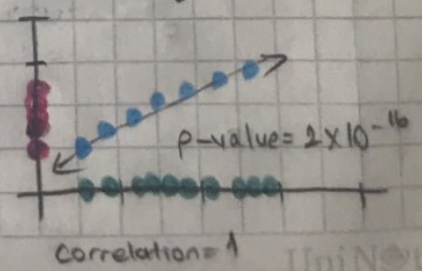
Korelasyon için p-value bize rastgele atılmış noktaların benzer şekilde güçlü bir ilişki veya daha güçlü bir ilişki ile sonuçlanma olasılığını söyler. Böylece, p-value ne kadar küçük olursa, doğru ile yaptığımız tahminlere o kadar güveniriz.



Bu case'te, p-value çok küçük: 2.2×10^{-16} . Bu rastgele verilerin benzer şekilde güçlü veya daha güçlü bir ilişki yaratma olasılığının oldukça küçük olduğu anlamına gelir.



Correlation = 1

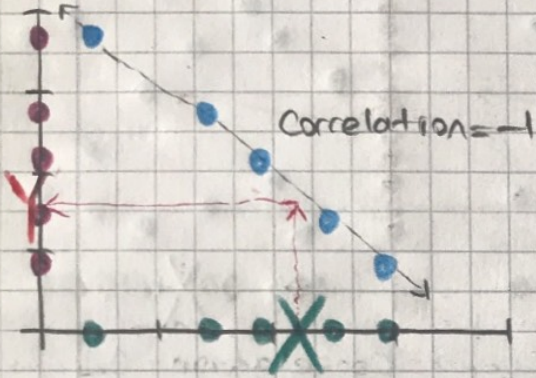


Correlation = 1

Veri arttikca p-value azalir ve bu, bizim korelasyon degerine olan guvenimizi artiriyor.

Negatif egimli bir dogru her veri noktasinin merkezinden gecbildiginde korelasyon = -1'dir.

* Bu bir dogru tum veri noktalarindan gecbildiginde, korelasyon verileri guclu bir iliski oldugu anlamina gelir ve eger birisi bize bir X' in bir deger verirse, bize Y' in degeri de bir aralikta bir deger tahmin edebiliriz.



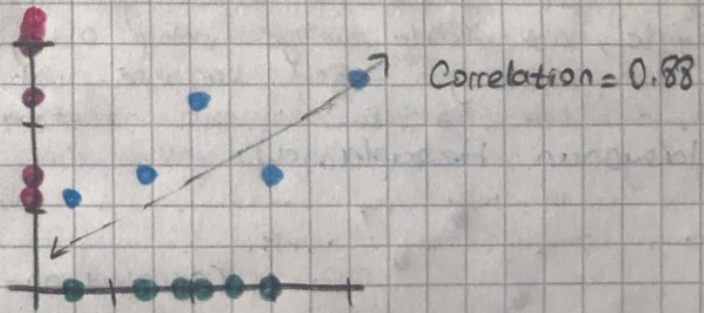
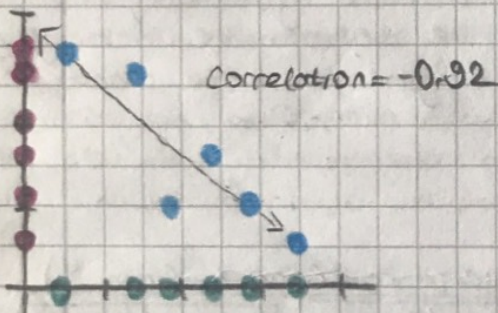
Ayni sekilde, p-value ile nercelelendirmede bu tahmine olan guvenimiz, elimizde ne kadar veri olduguna baglidir.

Gok fazla verim olsaydi, tahminimize cok guverebiliriz cunku p-value cok kucuk olurdu. (p-value = 2×10^{-16})

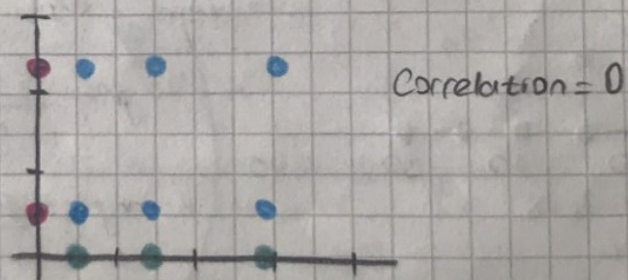
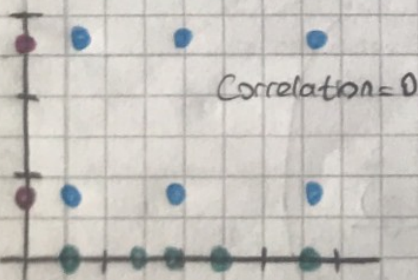
Ne kadar az veriye sahipsek, tahminimize o kadar az guvenimiz olur cunku p-value buyur. (p-value = 1)

Ayni sekilde negatif egimli bir dogru tum verilerden gecbildiginde korelasyon = -1 olur. Bu egim buyuk olsa da kucuk olsa da korelasyonun -1'e erit olabilecegi anlamina geliyor.

* Simdiye kadar, eger line'in egrisi pozitifse en guclu iliskinin correlation = 1 olacagini, eger line'in egrisi negatifse en guclu iliskinin correlation = -1 olacagini gorduk. Her iki durumda da eger dogru bota data larin uzerinden gecmiyorsa korelasyon degerleri 0'a yaklasiyor.

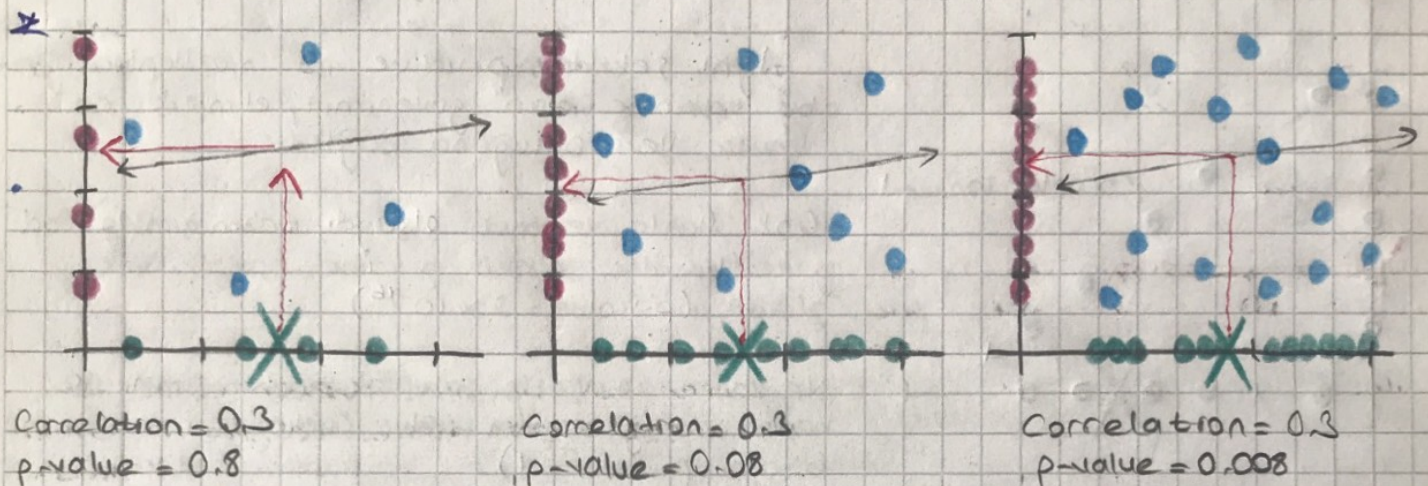


Veriler arasinda bir iliski olmadiginda correlation = 0 oluyor.



Correlation = 0 olduğunda x eksenindeki bir değer, y ekseninde hangi değeri beklememi gerektiği hakkında bize hiçbir şey söylemez. Çünkü x eksenindeki bir değere karşılık y ekseninde bir şeyi söylemek diğerini seçmemi için bir sebep yok.

Korrelasyon değeri 0 olmadığı sürece, hala çıkarım yapmak için lineer kullanabiliriz ancak korrelasyon değerleri -1 veya 1'e yaklaştıkça tahminlerimiz daha rafine (kırıktan uzak) hale gelir. Daha önce olduğu gibi, çıkarımlarımıza olan güvenimiz ver sayısına ve p-value'ya bağlıdır.



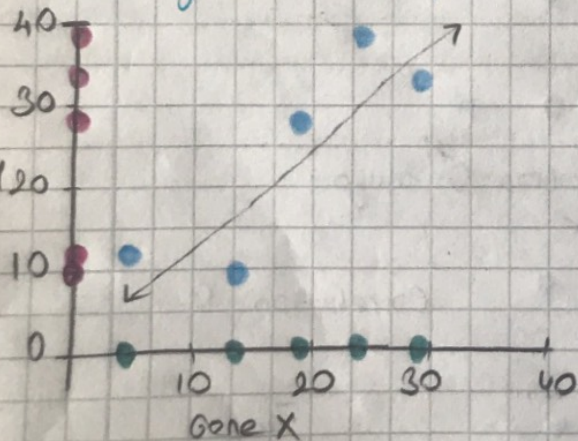
Verimiz az olduğu için bu trende çok az güveniyoruz.

Daha fazla verimiz oldu. Bu grafikte çok daha fazla verimiz olduğu için bu trende daha fazla güveniyoruz.

Bu grafikte çok daha fazla verimiz olduğu için trende çok daha fazla güveniyoruz.

Tüm örneklerde correlation = 0.3. Bu case'te, sample size'ni artırmak korrelasyonu artırmadı. Bu da veri eklemenin tahminimizi iyileştirmediği anlamına geliyor. Tek yaptığı tahmine olan güvenimizi artırmaktır. Bu nedenle, tahminlerimiz her üç durumda da muhtemelen kötü olacaktır. Ancak, en çok sorunlu verilerden gelen kötü tahmine güveneceğiz. Başka bir deyişle, çok fazla veriye sahip olduğunuz ve tahmininiz çok güvenilirdir, o zaman korrelasyon değeri küçükse tahminimiz yine de kötü olacaktır.

Korrelasyonun Hesaplanması

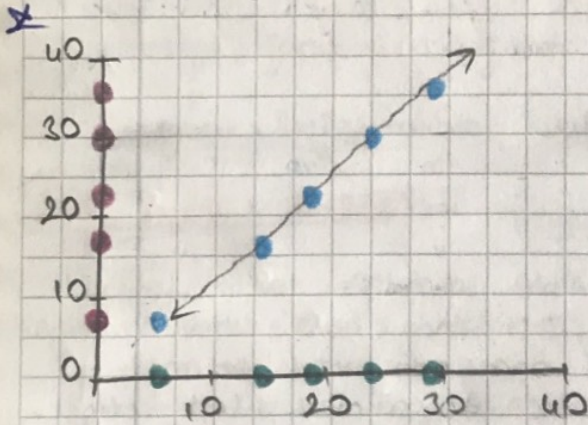


$$\text{Correlation} = \frac{\text{Covariance}(\text{Gene X}, \text{Gene Y})}{\sqrt{\text{Variance}(\text{Gene X})} \sqrt{\text{Variance}(\text{Gene Y})}}$$

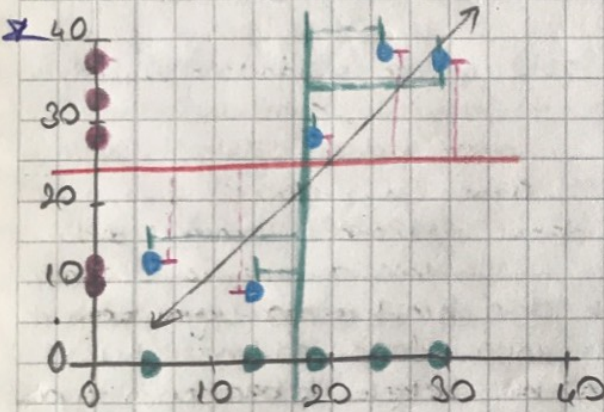
$$\text{Varyans} = \frac{\sum (x - \bar{x})^2}{n-1}$$

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

Paydaki covariance değeri -0.1 ile 0.1 arasında herhangi bir değer olabilir. Kısacası bu değer ilişkinin pozitif mi negatif mi olduğuna, verinin ortalama etrafında ne kadar uçukluğa dağıldığına ve verinin scale'ine bağlıdır. Böylece, korelasyonu hesapladığımızda, payda, kovaryansı -1 'den 1 'e kadar bir sayı olarak şekilde sınırlanır. Diğer bir deyişle, payda, verilerin ölçeğinin korelasyon değerini etkilememesini sağlar ve bu, korelasyonların yorumlanmasını çok daha kolay hale getirir.



Verilerin tümü pozitif veya negatif eğimli düz bir eğriye düştüğünde, kovaryans ve varyans terimlerinin kareköklerinin oranını aykır ve böyle eğime bağlı olarak bir 1 veya -1 verir.



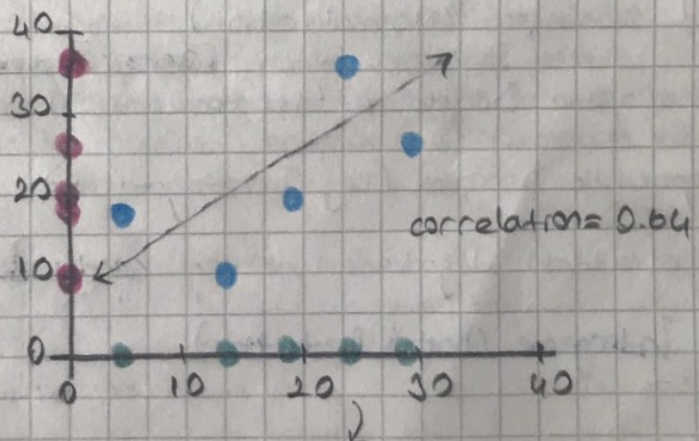
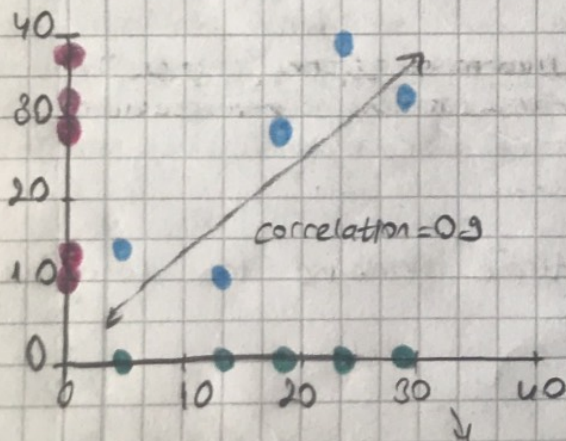
Veriler pozitif veya negatif eğimli düz bir eğriye düşmedikçe, kovaryans verilerdeki varyansın daha azını açıklar ve korelasyon 0 'a daha yakındır.

Bu veri için hesaplanacak olursa;

$$\text{Correlation} = \frac{11.6}{\sqrt{101.8} \sqrt{160.3}} = 0.9$$

$\text{correlation} = 0.9 \rightarrow$ bu ilişkiye olan güvenimi p -value ile ölçebiliriz. p -value ne kadar küçük olursa, yaptığımız tahminlere o kadar güvenebiliriz. Bu case'te $p\text{-value} = 0.03$. Bu, rastgele verilerin benzer şekilde veya daha güçlü bir ilişki üretme olasılığının $\%3$ olduğu anlamına gelir.

Korelasyon değerlerinin yorumlanması, kovaryans değerlerinden çok daha kolay olsa da, yorumlanmaları hala çok kolay değildir.



Örneğin, $\text{correlation} = 0.9$ olduğu bu ilişkinin, $\text{correlation} = 0.64$ olduğu bu ilişkiden tahmin yapmakten 1.5 kat daha iyi olduğu açık değildir. 2 bu sorunu çözüyor. UniNote