

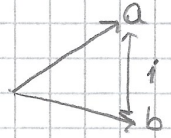
## Neden Normalization yapıyorum?

- \* Veriler arasında farklılığın çok fazla olduğu durumlarda verileri tek bir düzende ele alabilmek için
- \* Ayrıca, farklı ölçekleme sisteminde bulunan verilerin birbirini ile karşılaştırılabilmesi için.

- Complex sayılarla daha rahat çalışmak için.
- 0-1 aralığı için.

$$X = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \rightarrow \begin{matrix} 1500, 75000 \\ 20, 60 \end{matrix} \rightarrow \begin{matrix} 1500'den \\ 75000'e \end{matrix} \text{ stajyerden CEO'ya maaş}$$

$$KNN \quad d = \sqrt{(x_1^a - x_1^b)^2 + \dots + (x_n^a - x_n^b)^2}$$



## \* → Why doesn't we need normalization in iris data?

iris data da



genişlik ve boy bilgileri içerir ve bu bilgiler hep cm cinsi old. için burada normaliz. gerek yoktu. Bu dataset norm. olmadan kurtarır. Fakat diğer datasetlerde önce normal. yapmalısın.

Bazı datasetlerde verilerin birbirleri ile birimleri tutmayabilir. veya aralıkları çok yüksek olabilir. Mesela;

Yaş	Ücret
20	500
↓	↓
40	800

→ burada knn distance kullanırken

$$\sqrt{(800 - 500)^2 + (40 - 20)^2} \dots$$

bu çok büyük çıkacağı için diğer  $x_n$ 'leri önemsiz gibi gösterir

\* Hepsini dikkate alıp düzgünce çözümleyebilmek & gerçekçi olsun diye normalizasyon yapmalısın.

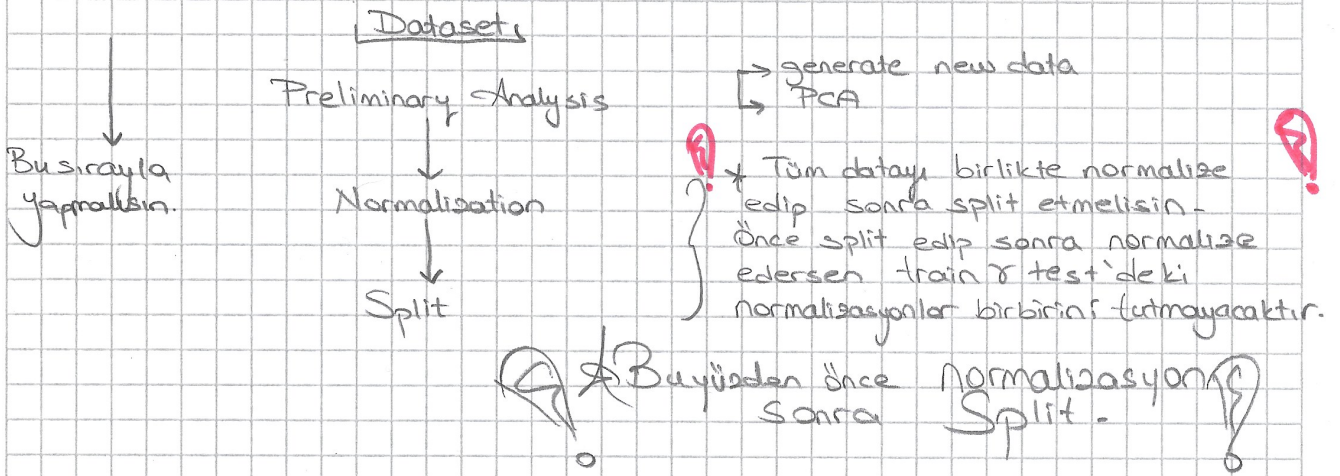
$$X = \frac{x - x_{\min}}{x_{\max} - x_{\min}}$$



→ [0,1] → 0-1 arasında olması için normalize yaparız.

$$\text{Standard scaler} = X = \frac{x - \mu}{\sigma}$$





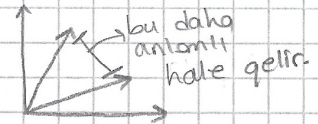
## Normalization, KNN, DT & SVM

### KNN'de normalization:

$x_1$ 'in ölçeğindeki küçük bir değişiklik  $x_2$ 'deki değişimi yok sayıyorsa yani meselim 10.000 iken 11.000'e çıktığında değişim  $(1000)^2$  iken Emağa göre küçük bir ölçek fakat sayısal olarak  $(bin)^2$  yüksek oldu için) yaş 20 → 60 değişim  $(40)^2$  olduğunda yaşı değişimi kendi skalasında çok olmasına  $(1000)^2$ 'nin yanında hiçbir önemi kalmaz. Bu yüzden normalize ediyoruz -

Yaş:	20	40	60
Normalize Yaş:	0	0.5	1

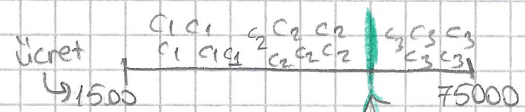
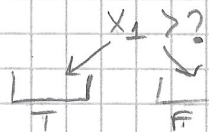
böyle yapınca



### Decision Tree'de normalization: min-max, standart dene hiçbir şey fark etmeyerek

Çünkü DT şöyle çalışır:

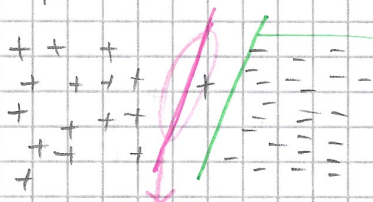
$C_1: 50$   
 $C_2: 50$   
 $C_3: 50$



\* DT ben bunu nereden ayırırsam en az hetero olur diyor. Buraya bir yer koysam ne kadar hetero olur diyor. Yani 0'dan 1'e gitmiş, 1500'den 75000'e gitmiş fark etmiyor. O yine dataset'i ortadan bir yerden bölüyor.

### Support Vector Machine'de (SVM) normalization:

$C = \text{margin}$



\* Ben buna baktığımda hard margin kullanırsam general-  
uak olurum. Soft margin kullanırsan sonuç daha smooth olur.

Benim samplem böyle ama  
benim görmediğim yerlerde belkide  
+, - karışık - var.

! C önemli !

$C = 1$  old. daha loose (?) öğrenme yapıyorsun  
özellikle birbirine geçmişi datalarda