

233. Kime Göre Neye Göre Aykırı Gözlem

Aykırı gözlemler veride genel eğilimin dışında çıkan gözlemler olarak tanımlanmıştır. Peki veri setinin genel eğiliminin dışında çıkmayı nasıl tanımlarız?

1. Sektör Bilgisi

Bazı senaryolarda aykırı gözlemleri nümerik yöntemlerle tanımlayamayız da çalışıyor olduğumuzu sektörün dinamikleriyle göre neyin aykırı ya da neyin normal olduğu durumunu tanımlayabiliriz.

Örneğin; bir ev fiyat tahmin modelinde 1000 metrekarelik evleri modellemeye almamak. Çünkü bizim makine öğrenmesi çalışmamızda amaçımız genelleşebilir, veri setinin içerisindeki yapıları tanımlı etmek kabiliyeti yüksek, yanlış modeller oluşturmak Türkiye için derinleşmiş bir makine belki evlerin %99'unun 2+1, 3+1 veya daha küçük olduğunu biliyoruz. Dolayısıyla belki %1 bile olmayan 1000 metrekarelik evleri bir tahmin modeli için çalışmaya dahil etmek, diğer daha düşük metre-karelerde yer alan evler için yapılacak olan genellemelere zarar verecektir.

Çünkü kullanıyor olduğumuz fonksiyonel yapılar sadece 1000 metrekarelik evlerle ilişkili yapılarla alakalı da bazı genelleştirilebilirlikler sağlamaya çalıştığımızdan dolayı, bu durumda veri setinin içerisindeki diğer yapıları kaçırabiliyoruz.

2. Standart Sapma Yaklaşımı

Bir değişkenin ortalamasının üzerine aynı değişkenin standart sapması hesaplanarak eklenir. 1, 2 ya da 3 standart sapma değeri ortalamaya üzerine eklenerek ortaya çıkan bu değer eşik değer olarak düşünülür ve bu değerden yukarıda ya da aşağıda olan değerler aykırı değer olarak tanımlanır.

Diyecek ki elimizde araç fiyatları var. Bunların ortalaması ve standart sapma bilgisi var elimizde. Bunları kullanarak eşik değeri şu şekilde belirleyebiliriz.

$$\text{Eşik Değer} = \text{Ortalama} + 1 \times \text{Standart Sapma}$$

$$\text{Eşik Değer} = \text{Ortalama} + 2 \times \text{Standart Sapma}$$

$$\text{Eşik Değer} = \text{Ortalama} + 3 \times \text{Standart Sapma}$$

Hongrını kullanacağımız problemleri probleme düşebiliriz.

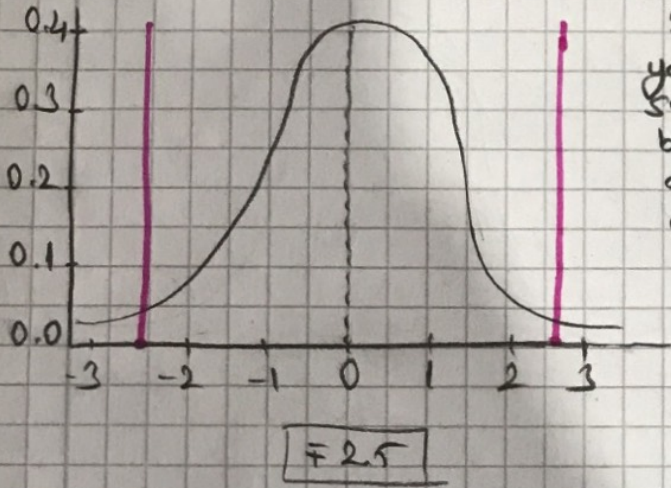
Örn; araç fiyat ortalaması 100.000 ve Standart Sapma = 20.000 ise bu durumda şu şekilde belirleyebiliriz. Benim verimin merkezi eğilimini ben buldum 100.000. Üzerine bu değişkenin dağılımını (ortalama etrafındaki yapılamamı) ifade eden standart sapmadan 2 tane ekledim.

$$\begin{aligned} \text{Eşik Değer} &= 100.000 + 2 \times 20.000 \\ &= 140.000 \end{aligned}$$

3.2 - Skoru Yaklaşımı

Eğer elimizdeki değeri standartlaştırılıp buna göre bir muamele yaparsak bu durumda z-skoru kullanılır.

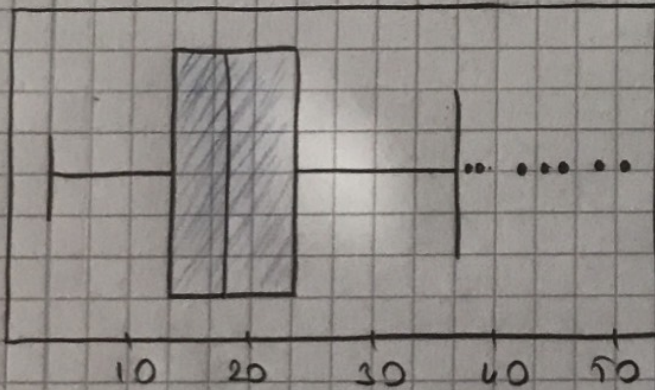
Standart sapma yöntemine benzer şekilde yapılır. Değerler, standart normal dağılıma uyarlanır, yani standartlaştırılır. Sonrasında -örneğin- dağılımın sağından ve solundan ± 2.5 değerine göre bir eşik değeri konulur ve bu değerin uzerinde ya da altında olan değerler aykırı değer olarak işaretlenir.



Buradaki yaklaşımda, değerin genel yapının uzerine, değerin genel yapının dağılımını da ifade eden $+$, $-$ bir değer koyarak veri setinin genel dağılımının dışındaki yapılara erişilme sağlanır.

4. Boxplot (interquartile range - IQR) Yöntemi

En sık kullanılan yöntemlerden birisidir. Değerlerin değerleri küçükten büyüğe sıralanır. Gerekliliklerine (yadeliiklerine) yani Q_1 , Q_3 değerlerine karşılık gelen değerler uzerinden bir eşik değeri hesaplanır ve bu eşik değere göre aykırı değer tanımlı yapılır.



Bu grafik minimum değeri, maksimum değeri, 3. çeyrekteki değeri, 1. çeyrekteki değeri ve medyana göre oluşturuluyordu.

Burada eşik değerleri hesaplamak için;

$$IQR = 1.5 \times (Q_3 - Q_1)$$

değeri hesaplanır. Bu değeri alt eşik değeri ve üst eşik değeri hesaplanırken aşağıdaki şekilde kullanılır.

$$\text{Alt Eşik Değeri} = Q_1 - IQR$$

$$\text{Üst Eşik Değeri} = Q_3 + IQR$$

Genellikle 1. veya 4. yaklaşım kullanılır.