

# Keşifsel Veri Analizi (EDA) - ML

🔧 Veri Temizleme (Data Cleaning)

## Keşifsel veri analizine giriş

### Keşifsel Veri Analizi (Exploratory Data Analysis) Nedir?

İstatistikte, keşifsel veri analizi, genellikle istatistiksel grafikler ve diğer veri görselleştirme yöntemlerini kullanarak temel özelliklerini özetlemek için veri kümelerini analiz etme yaklaşımıdır. İstatistiksel bir model kullanılabilir veya kullanılamaz.

### EDA neden kullanışlıdır?

- EDA veriler için ilk izlenim ve gözlem oluşturmamızı sağlar.
- Bu şekilde verilerin anlamlı olup olmadığını ya da daha fazla temizliğe ihtiyacı olup olmadığını anlayabiliriz.
- Verilerdeki kalıpları ve eğilimleri belirlememize yardımcı olur.

### EDA teknikleri

- **Özetleyici nitelikler** - Veri hakkında bilgilenmemizi sağlayacak özet verilerdir. Bunlar ortalama, medyan , maksimum değer , korelasyon vs.

- **Görselleştirme nitelikleri** - Histogramlar , Saçılım grafikleri, kutu grafikleri vs.

Veriyi incelemek ve tartışmak için pythonda **Pandas** kütüphanesini kullanacağız.

Veriyi görselleştirmek içinse **Matplotlib** ve **Seaborn** kütüphanelerini kullanacağız.,

Verinin bazı örneklerini karışık olarak görebilmek için aşağıdaki kod kullanılır.

```
örnekler = veri.sample(n=5 , replace = False)
print(örnekler.iloc[:,-3:])
```

## Veri görselleştirme kütüphaneleri

- Matplotlib - Matplotlib temel grafikleri çizdirmemizi ve veri hakkında bilgi edinmemizi sağlayan grafik kütüphanesidir.
- Pandas (via Matplotlib) - Pandas kütüphanesinde veri görselleştirmede kullanılabilir ancak matplotlib kadar hızlı ve esnek değildir. Yinede gayet yeterlidir.
- Seaborn - Seaborn Kütüphanesi fantastik grafikler çizdirmek için kullanılabilir. Matplotlibe göre daha büyüleyici temalar kullanırlar. Matplotlibden dahil edilmiş bazı içerikler barındırır. Yani seabornda temel olarak matplotlibi kullanır.

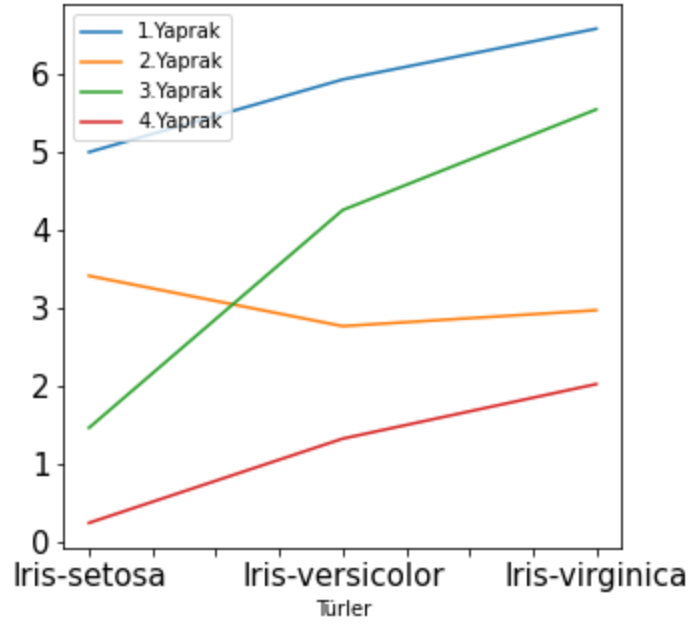
## Veri Görselleştirme

### 1. Gruplandırma

Verileri sınıfına göre gruplandırarak veriler hakkında bilgi sahibi olabilir. Burada iris veri kümesini ele alacağız.

```
veri.groupby("Türler").mean().plot(fontsize=15,figsize=(5,5))
```

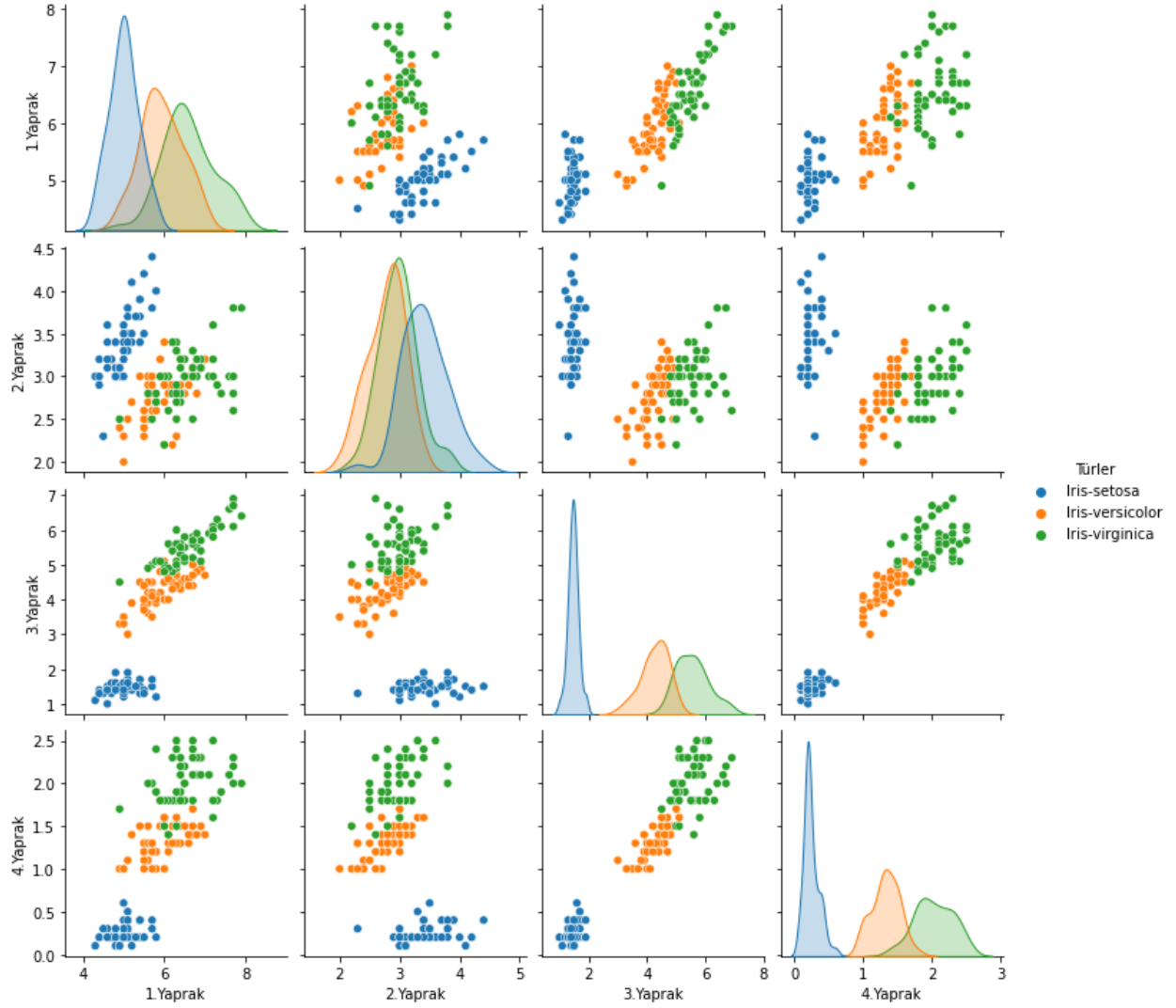
Yukarıda ki kod karşılığında çıktımız aşağıdaki şekilde olacaktır. Bu grafik hangi değişkenin hangi durumda nasıl bir davranış sergilediğini gösterir.



## 2. Çiftler Grafiği

Özelliklerin kendi aralarındaki bağlantıyı görmek için kullanabileceğimiz bir veri görselleştirme türüdür. Burada sınıflar esas alınarak özellikler arasındaki bağlantı gözlenir. Bunun için aşağıdaki python kodu kullanılır. Girdi olarak verinin kendisini vermek yeterlidir. renk tonu olarak sınıfların verilmesi, verinin görselleştirilmesi açısından çok daha iyi olacaktır

```
import seaborn as sns
sns.pairplot(veri, hue="Türler")
```

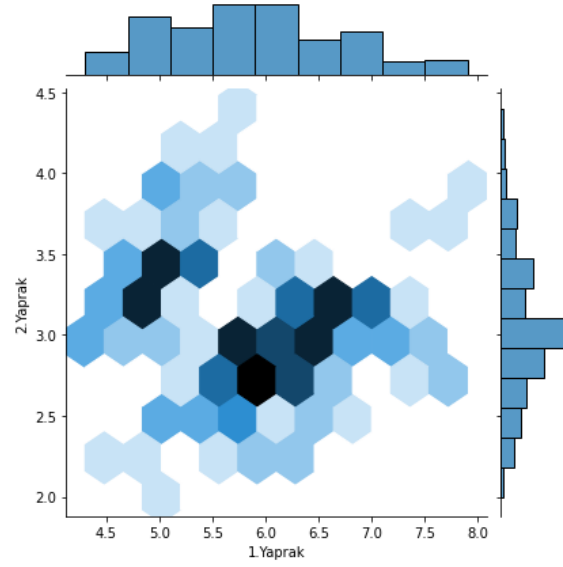


### 3. Hekzagonal (Altıgen) Grafik oluřturma

İki veri arasındaki baęlantıyı g rebilmek i in kullanılabilecek grafik t rlerinden biridir. kind parametresinin alacaęı deęere g re grafięin řekli deęiřir.  rneęin kind= "kde" se erseniz yoęunluęa g re bir izohips grafięi  ıktısı alınır.

```
sns.jointplot(veri["1.Yaprak"], veri["2.Yaprak"], kind = "hex")
```

Yukarıdaki kod yazıldıęında  ıktımız ařaęıdaki gibi olacaktır. Yoęunluęun y ksek olduęu yerlerde renkte koyulařma olur.

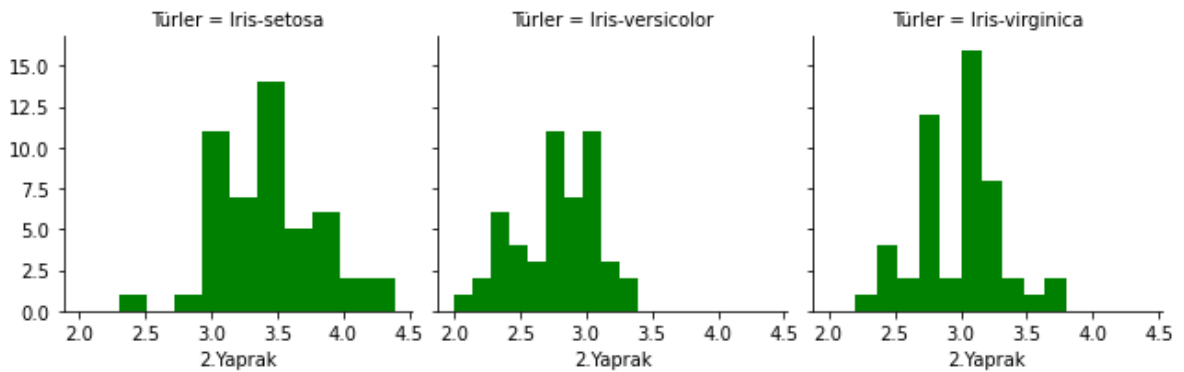


#### 4. Izgara methodu

Eğer her bir kolonun sınıf üzerindeki etkisini görmek istiyorsak bunun için Facet\_grid kullanılır.

```
grafik = sns.FacetGrid(veri,col="Türler")
grafik.map(plt.hist,"2.Yaprak",color="green")
```

Bu kodun çıktısı aşağıdaki şekilde olacaktır.

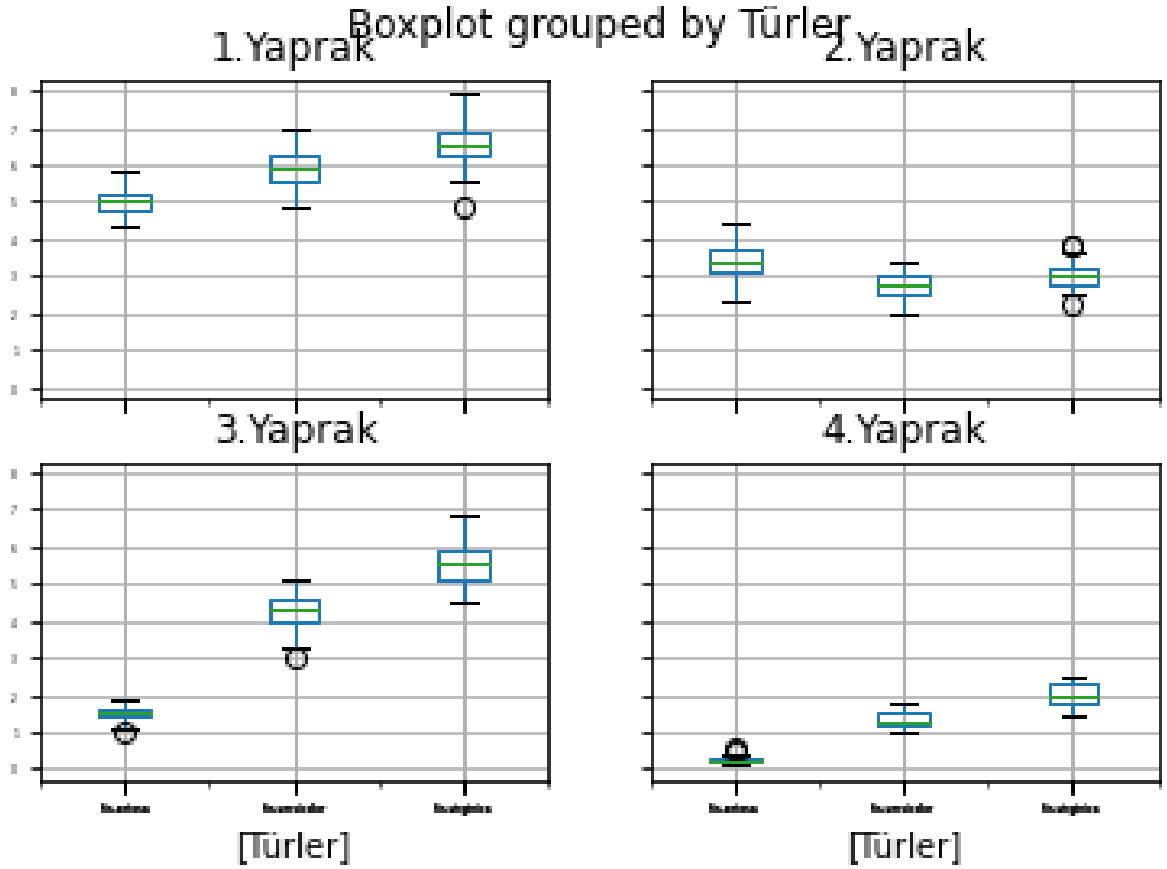


#### 5. Kutu Grafiği

Verileri istediğimiz referansa göre kutu biçiminde göstermemizi sağlayan grafik methodudur. Aşağıdaki kod ile kullanılabilir.

```
veri.boxplot(by="Türler")
```

Bu kodun çıktısı aşağıdaki şekilde olacaktır.



Değişken dönüşümleri

Mert Aydoğan Gazi University Department of ME