



## ***MAKİNE ÖĞRENMESİ FİNAL PROJESİ***

### **HAZIRLAYAN**

**ADI SOYADI** : FURKAN DOĞAN

**ÖĞRENCİ NUMARASI** : 180303026

**TESLİM TARİHİ** : 06/01/2023

**DERS ADI** : MAKİNE ÖĞRENMESİ

**DERS YÜRÜTÜCÜSÜ** : DR. SAEİD AGAHİAN

## **BÖLÜM – 1 : VERİ SETİNE GENEL BAKIŞ**

**a. Problemin Tanımı**

**b. Sınıf ve Örnek Sayısı**

## **BÖLÜM - 2 : ALGORİTMA HAKKINDA TEORİK BİLGİ**

## **BÖLÜM – 3 : ÇAPRAZ GEÇERLİLİK (CROSS-VALIDATION) DEĞERLENDİRME ÖLÇEĞİ**

## **BÖLÜM – 4 : CONFUSION MATRIX İLE DEĞERLENDİRME**

## 1 - VERİ SETİNE GENEL BAKIŞ

Ses veri setinden Random Forest algoritmasını kullanarak bir Cinsiyet Tanıma sınıflandırıcısı oluşturacağız. Fikir, sesin ve konuşmanın akustik özelliklerine dayalı olarak bir sesi erkek veya dişi olarak tanımlamaktır. Veri seti, erkek ve kadın konuşmacılardan toplanan 3.168 kayıtlı ses örneğinden oluşmaktadır. Ses örnekleri, 0hz-280hz analiz frekans aralığı ile Seewave ve tuneR paketleri kullanılarak R'de akustik analiz ile önceden işlenmiştir. **Veri seti Kaggle' den alınmıştır.**

### Problemin Tanımı

Amaç, bir Karar ağacı ve Rastgele Orman Ağacı sınıflandırıcısı oluşturmak ve her iki modelin doğruluğunu karşılaştırmaktır. 'label' sütunu, kişinin cinsiyetini tahmin etmemiz gereken hedef değişkendir.

### Sınıf ve Örnek Sayısı

- **meanfreq**: ortalama frekans (kHz cinsinden )
- **sd**: frekansın standart sapması
- **median**: medyan frekans (kHz cinsinden)
- **Q25**: birinci nicelik (kHz cinsinden)
- **Q75**: üçüncü nicelik (kHz cinsinden)
- **IQR**: çeyrekler arası aralık (kHz cinsinden)
- **skew**: çarpıklık
- **kurt**: basıklık
- **sp.ent**: spektral entropi
- **sfm**: spektral düzlük
- **mode**: mod frekansı
- **centroid**: frekans merkezi
- **peakf**: tepe frekansı (en yüksek enerjiye sahip frekans)
- **meanfun**: akustik bir sinyal boyunca ölçülen temel frekansın ortalaması
- **minfun**: akustik bir sinyal boyunca ölçülen minimum temel frekans
- **maxfun**: akustik bir sinyal boyunca ölçülen maksimum temel frekans
- **meandom**: bir akustik sinyal boyunca ölçülen baskın frekansın ortalaması
- **mindom**: bir akustik sinyal boyunca ölçülen minimum baskın frekans
- **maxdom**: bir akustik sinyal boyunca ölçülen maksimum baskın frekans
- **dfrange**: bir akustik sinyal boyunca ölçülen baskın frekans aralığı
- **modindx**: frekans aralığına bölünen temel frekansların bitişik ölçümleri arasındaki birikmiş mutlak fark olarak hesaplanan modülasyon indeksi
- **label**: bay veya bayan

21 tane sınıf bulunmaktadır.

## 2 - ALGORİTMA HAKKINDA TEORİK BİLGİ

### Decision Tree Classification

Ağaç tabanlı öğrenme algoritmaları, en çok kullanılan ve denetimli öğrenme yöntemlerinden biri olarak düşünülmektedir. Ağaç tabanlı yöntemler, yüksek doğruluk, kararlılık ve yorumlanma kolaylığına sahiptir. **Decision Tree** yapısında, ağaç üzerindeki her bir karar noktasında veri kümesinin o anki dağılımına göre bir "ölçüt" kullanılır. Bu ölçüt, ağaç yapısının ne şekilde büyüyeceğini belirler. Bu ölçütlerden ikisi, **entropy** ve **gini** indeksleridir.

**Entropy**: Veri kümesindeki örneklerin sınıf dağılımına göre bir ölçüttür. Eğer veri kümesinde sınıflar eşit dağılım gösteriyorsa, **entropy** değeri yüksek olur. Eğer veri kümesinde bir sınıf aşırı derecede öne çıkıyorsa, **entropy** değeri düşük olur.

**Gini** indeksi: **Gini** indeksi de veri kümesindeki örneklerin sınıf dağılımına göre bir ölçüttür. Eğer veri kümesinde sınıflar eşit dağılım gösteriyorsa, **gini** indeksi değeri yüksek olur. Eğer veri kümesinde bir sınıf aşırı derecede öne çıkıyorsa, **gini** indeksi değeri düşük olur. Hangisi kullanılacağı, veri kümesine göre değişebilir. Örneğin, eğer veri kümesinde sınıflar eşit dağılım gösteriyorsa, **entropy** değeri yüksek olacağı için **entropy** tercih edilebilir. Eğer veri kümesinde bir sınıf aşırı derecede öne çıkıyorsa, **gini** indeksi değeri düşük olacağı için **gini** indeksi tercih edilebilir.

**Gini indeksi(dizini) veya Gini katsayısı** : İtalyan istatistikçi Corrado Gini tarafından 1912’de geliştirilen istatistiksel bir ölçüdür. Katsayı, 0 (%0) ile 1 (%100) aralığındadır; 0, mükemmel eşitliği temsil eder ve 1, mükemmel eşitsizliği temsil eder.

**Entropy** : rasgele bir değişkenin belirsizliğinin ölçüsüdür, örneklerin keyfi bir koleksiyonunun saf olmayanlığını karakterize eder. Entropi ne kadar yüksek olursa elde edilen bilgi de o kadar fazla olur. Sezgisel olarak, belirli bir olayın öngörülebilirliğinden bahseder.

Entropi, tipik olarak, eğitim örneklerini daha küçük alt gruplara bölmek için bir karar ağacında bir düğümü kullandığımızda değişir. Bilgi kazancı, entropideki bu değişimin bir ölçüsüdür.

## **Random Forest Classification**

Rastgele ormanlar veya rastgele karar ormanları, sınıflandırma, regresyon ve diğer görevler için, eğitim aşamasında çok sayıda karar ağacı oluşturarak problemin tipine göre sınıf veya sayı tahmini yapan bir toplu öğrenme yöntemidir. Random forest, birden fazla karar ağacını kullanarak daha uyumlu modeller üreterek daha isabetli sınıflandırma yapmaya çalışır. Bu yüzden öncelikle DT Sınıflandırmasını ardından RF sınıflandırmasını kullandım. Çalışma prensibi ise ;

[ DT = Decision Trees - RF = Random Forest ] adından da anlaşılacağı gibi, bir orman oluşturur ve bunu bir şekilde rastgele yapar. Kurduğu “orman”, çoğu zaman “bagging” yöntemiyle eğitilen karar ağaçları topluluğudur. Bagging yönteminin genel fikri, öğrenme modellerinin bir kombinasyonunun genel sonucu arttırmasıdır.

## **K-Nearest Neighbors(KNN)**

Diğer Denetimli Öğrenme algoritmalarının aksine, eğitim aşamasına sahip değildir. Eğitim ve test hemen hemen aynı şeydir. Tembel bir öğrenme türüdür. Bu nedenle, kNN, geniş veri setini işlemek için gereken algoritma olarak ideal bir aday değildir. KNN ile temelde yeni noktaya en yakın noktalar aranır. K, bilinmeyen noktanın en yakın komşularının miktarını temsil eder. Sonuçları tahmin etmek için algoritmanın k miktarını (biz k=10 kabul edeceğiz) seçeriz. KNN algoritmasının bir özelliği, verilerin yerel yapısına duyarlı olmasıdır.

## **3 - ÇAPRAZ GEÇERLİLİK (CROSS-VALIDATION) DEĞERLENDİRME ÖLÇEĞİ**

Cross-validation, makine öğrenmesi modelinin görmediği veriler üzerindeki performansını mümkün olduğunca objektif ve doğru bir şekilde değerlendirmek için kullanılan istatistiksel bir yeniden

örnekleme(resampling) yöntemidir. İkinci bir kullanım alanı ise modelde hiper parametre optimizasyonu yapmaktır. K-Fold CV yöntemi üzerinden çalışma mantığı ise şudur :

1. Veri seti rastgele olacak şekilde karıştırılır. (opsiyonel)
2. Veri seti k gruba ayrılır.
3. Her grup için aşağıdaki adımlar uygulanır :
  - Seçilen grup validasyon seti olarak kullanılır.
  - Diğer tüm gruplar (k-1 grup) train seti olarak kullanılır.
  - Train seti kullanılarak model kurulur ve validasyon seti ile değerlendirilir.
  - Modelin değerlendirme puanı bir listede saklanır.

Şimdi projeye ait kısımları inceleyelim..

**Yorum satırlarında yapılan her şey açıklanmıştır.**

```
1 #Gerekli paketler yüklendi.  
2 !pip install opencv-python  
3 !pip install matplotlib-venn  
4 !apt-get -qq install -y libfluidsynth1  
5 !pip install -U scikit-learn  
6 !pip install scikit-learn  
7
```

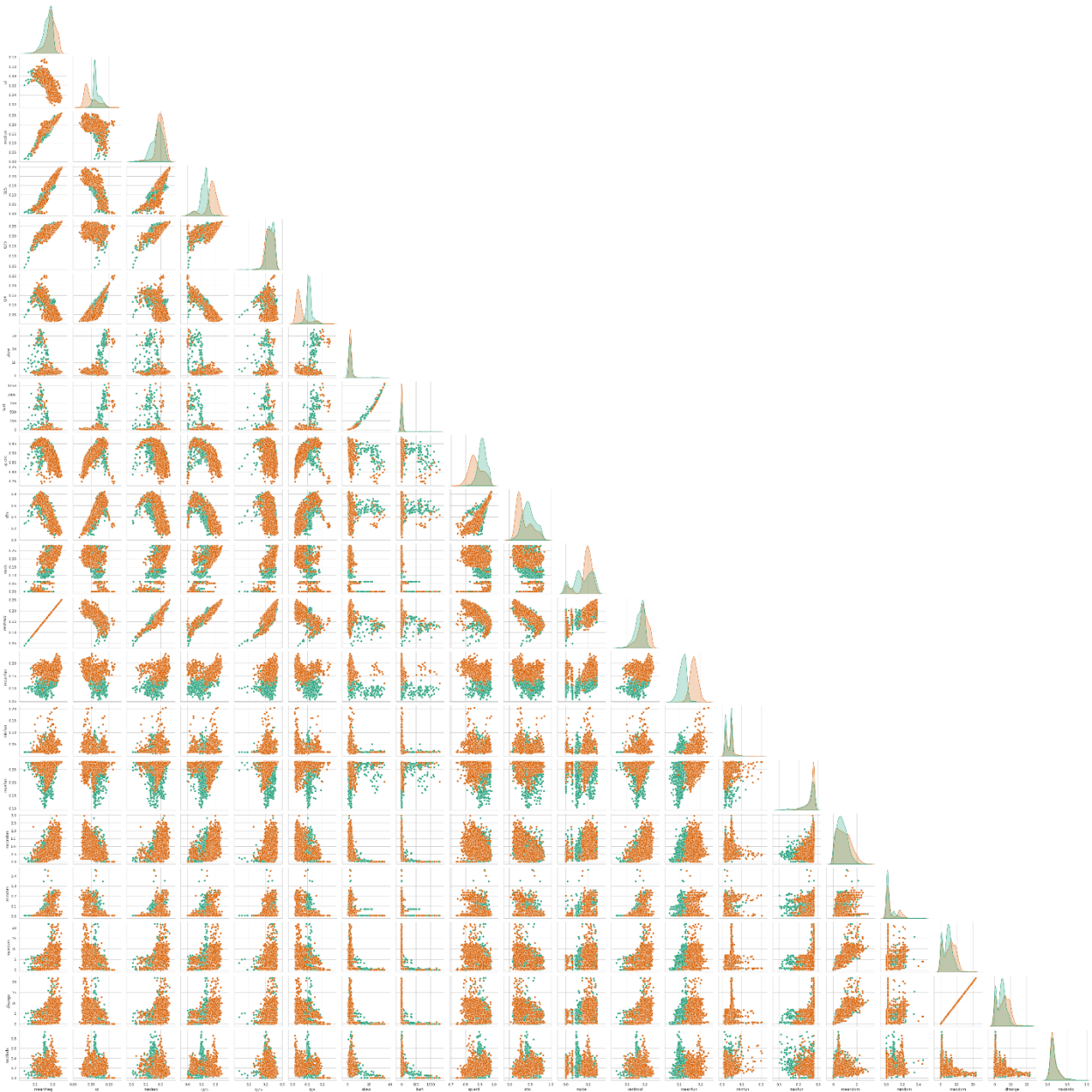
## Rapor/Ödev Başlığı: BİL 411 / Makine Öğrenmesi Final Projesi

### Hazırlayanın Adı Soyadı: Furkan Doğan

```
1 train.head().style.background_gradient(cmap='coolwarm') # Veri kümesinin ilk beş satırını göster ve değerleri renklerle belirt.
```

	meanfreq	sd	median	Q25	Q75	IQR	skew	kurt	sp.ent	sfm	mode	centroid	meanfun	minfun	maxfun	meandom	mindom	maxdom	dfrange	modindx	label
0	0.059781	0.064241	0.032027	0.015071	0.090193	0.075122	12.863462	274.402906	0.893369	0.491918	0.000000	0.059781	0.084279	0.015702	0.275862	0.007812	0.007812	0.007812	0.000000	0.000000	male
1	0.066009	0.067310	0.040229	0.019414	0.092666	0.073252	22.423285	634.613855	0.892193	0.513724	0.000000	0.066009	0.107937	0.015826	0.250000	0.009014	0.007812	0.054688	0.046875	0.052632	male
2	0.077316	0.083829	0.036718	0.008701	0.131908	0.123207	30.757155	1024.927705	0.846389	0.478905	0.000000	0.077316	0.098706	0.015656	0.271186	0.007990	0.007812	0.015625	0.007812	0.046512	male
3	0.151228	0.072111	0.158011	0.096582	0.207955	0.111374	1.232831	4.177296	0.963322	0.727232	0.083878	0.151228	0.088965	0.017798	0.250000	0.201497	0.007812	0.562500	0.554688	0.247119	male
4	0.135120	0.079146	0.124656	0.078720	0.206045	0.127325	1.101174	4.333713	0.971955	0.783568	0.104261	0.135120	0.106398	0.016931	0.266667	0.712812	0.007812	5.484375	5.476562	0.208274	male

Nitelikler arasındaki korelasyonu ve sınıf değişkenine göre dağılımı gösteren çoklu görsel.







## CONFUSION MATRIX İLE DEĞERLENDİRME

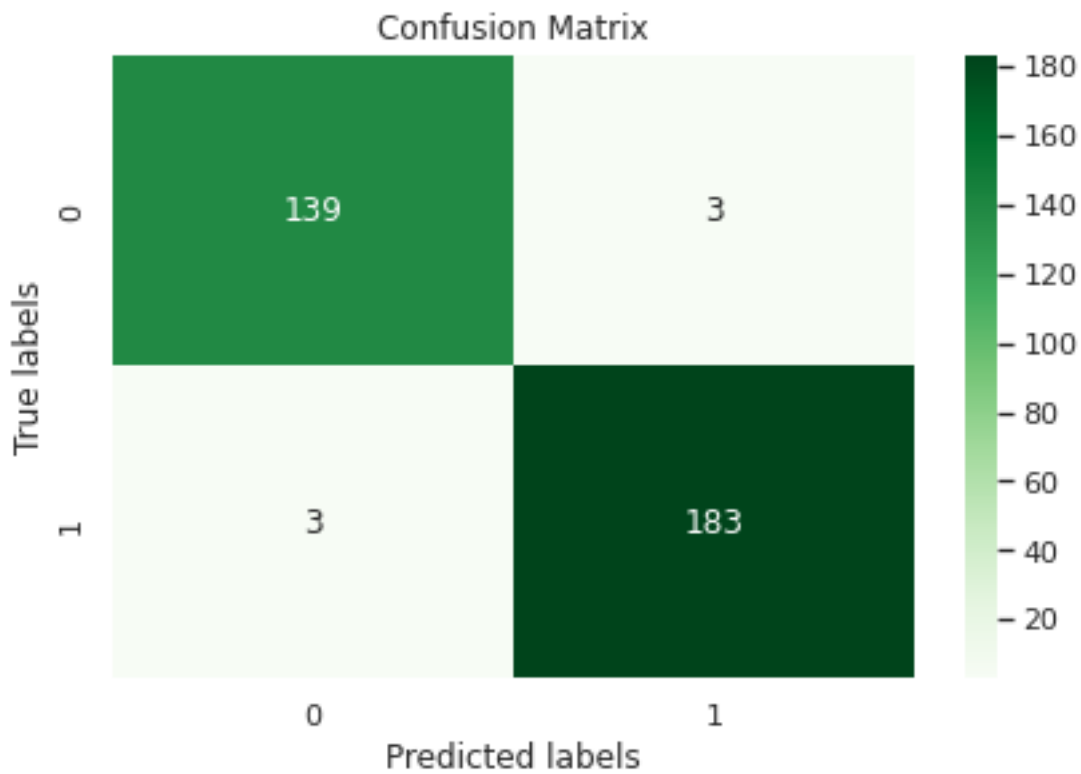
Makine öğrenmesinde kullanılan sınıflandırma modellerinin performansını değerlendirmek için hedef niteliğe ait tahminlerin ve gerçek değerlerin karşılaştırıldığı hata matrisinin görseli aşağıdadır. Her ne olursa olsun sınıflandırma tahminleri şu dört değerlendirmeden birine sahip olacaktır:

- Doğruya doğru demek (True Positive – TP) **DOĞRU**
- Yanlışta yanlış demek (True Negative – TN) **DOĞRU**
- Doğruya yanlış demek (Predicted Positive – FP) **YANLIŞ**
- Yanlışta doğru demek (Predicted Negative – FN) **YANLIŞ**

## RANDOMFORESTCLASSIFIER CONFUSION MATRIX

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 pred = clf_rf.predict(x_test) # modelden tahminler yap.
5 conf_matrix = confusion_matrix(y_test, pred) # confusion matrix oluştur.
6
7 #RandomForestClassifier confusion matrix'i grafik şeklinde gösterin
8 ax= plt.subplot()
9 # confusion matrix'i "Greens" renk paletiyle göster.
10 sns.heatmap(conf_matrix, annot=True, ax = ax, cmap='Greens', fmt='g')
11
12 # Eksen etiketlerini ayarla.
13 ax.set_xlabel('Predicted labels')
14 ax.set_ylabel('True labels')
15 ax.set_title('Confusion Matrix')
16 ax.xaxis.set_ticklabels(['0', '1'])
17 ax.yaxis.set_ticklabels(['0', '1'])
18
```

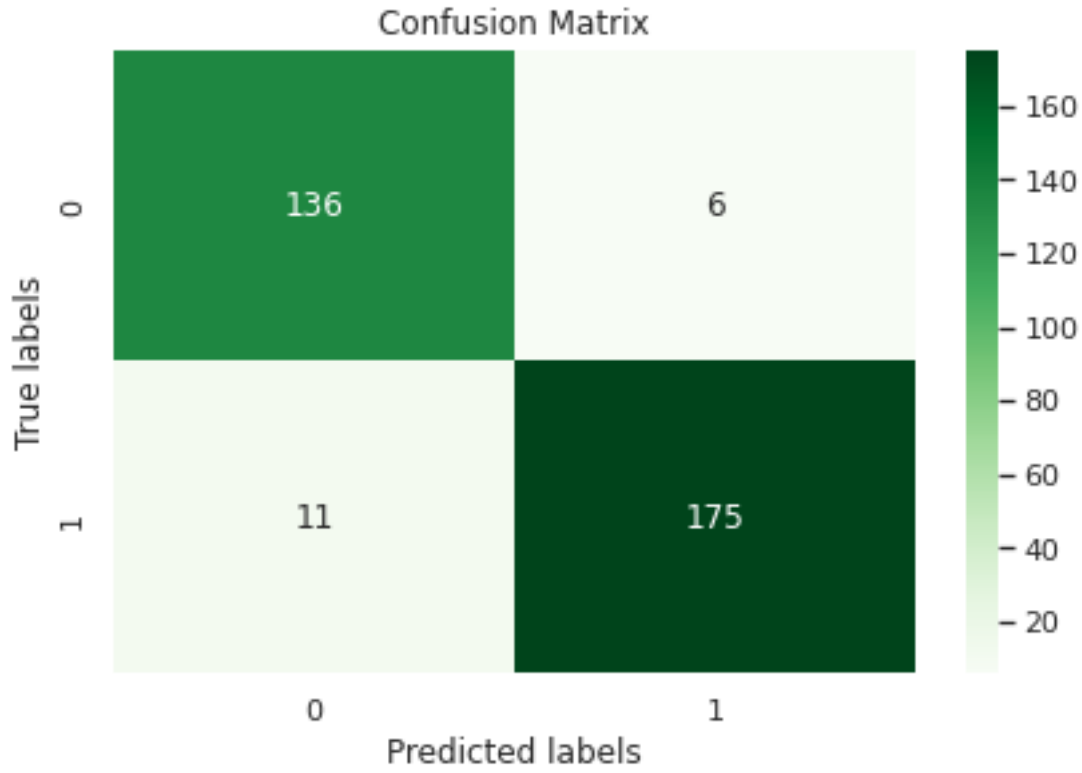
Yukarıdaki kodun çıktısı olarak aşağıdaki Confusion Matrix görüntüsü elde edilir.



## DECISIONTREECLASSIFIER CONFUSION MATRIX

```
1 import matplotlib.pyplot as plt
2 import seaborn as sns
3
4 pred = clf_dt.predict(x_test) # modelinizden tahminler yap.
5 conf_matrix = confusion_matrix(y_test, pred) # confusion matrix oluştur.
6
7 #DecisionTreeClassifier confusion matrix'i grafik şeklinde göster.
8 ax= plt.subplot()
9 # confusion matrix'i "Greens" renk paletiyle göster.
10 sns.heatmap(conf_matrix, annot=True, ax = ax, cmap='Greens', fmt='g')
11
12 # Eksen etiketlerini ayarla.
13 ax.set_xlabel('Predicted labels')
14 ax.set_ylabel('True labels')
15 ax.set_title('Confusion Matrix')
16 ax.xaxis.set_ticklabels(['0', '1'])
17 ax.yaxis.set_ticklabels(['0', '1'])
```

Yukarıdaki kodun çıktısı olarak aşağıdaki Confusion Matrix görüntüsü elde edilir.



## ACCURACY\_SCORE İLE DOĞRULUK ORANI

Yüksek bir doğruluk oranı, modelin etkili bir şekilde sınıflandırma yaptığı anlamına gelebilir.

```
1 from sklearn import tree
2
3
4 # Veri çerçevesinden etiket değişkenini düşürün ve değişkenleri standartlaştır.
5 scaler=StandardScaler()
6 scaled_df=scaler.fit_transform(temp_df.drop('label',axis=1))
7 X=scaled_df
8 Y=df['label']
9
10 # Eğitim ve test verilerini ayırın
11 x_train,x_test,y_train,y_test=train_test_split(X,Y,test_size=0.20,random_state=42)
12
13 #Bu kod parçasığında, veri kümesinden bir karar ağacı ve rastgele ormanı kullanarak sınıflandırma modelleri eğitilmektedir.
14 # Karar ağacı modelini oluşturun ve eğit.
15 clf_dt=DecisionTreeClassifier()
16 clf_dt.fit(x_train,y_train)
17
18 # Karar ağacı modelini görsel olarak çizdir.
19 plt.figure(figsize=(30,18))
20 tree.plot_tree(clf_dt,filled=True)
21 plt.show()
22
23 # Karar ağacı modelinin test verisi üzerinden tahminler yapın ve doğruluk skoru hesapla.
24 pred=clf_dt.predict(x_test)
25 print("Karar ağacı modeli doğruluk oranı: ",accuracy_score(pred,y_test))
26
27 # Rastgele orman modelini oluşturun ve eğit.
28 clf_rf=RandomForestClassifier()
29 clf_rf.fit(x_train,y_train)
30
31 # Rastgele orman modelinin test verisi üzerinden tahminler yapın ve doğruluk skoru hesapla.
32 pred=clf_rf.predict(x_test)
33 print("Rastgele orman modeli doğruluk oranı: ", accuracy_score(pred,y_test))
```

```
Karar ağacı modeli doğruluk oranı:  0.948170731707317
Rastgele orman modeli doğruluk oranı:  0.9817073170731707
```

"GridSearchCV()" fonksiyonu, modelinizin performansını optimize etmek için verilen parametrelerin değerlerini değiştirerek modeli eğitir ve test eder.

```
1 # En iyi skor ve parametreleri yazdır.  
2 print("Best score : ",CV_rfc.best_score_)  
3 print("Best Parameters : ",CV_rfc.best_params_)  
4 # Test verisindeki tahminlerin doğruluk skorunu yazdır.  
5 print("Precision Score : ", precision_score(CV_rfc.predict(x_test),y_test))
```

Best score : 0.9816472171039162

Best Parameters : {'criterion': 'entropy', 'max\_depth': 7, 'max\_features': 'auto', 'n\_estimators': 200}

Precision Score : 0.9838709677419355