

Makine Öğrenmesi Projesi Raporu

Proje Adı: Diyabet Veri Seti Üzerinde Makine Öğrenmesi Modelleri ile Tahminleme

Bu proje, diyabet hastalığı tahminine yönelik olarak makine öğrenmesi modellerinin karşılaştırılması ve en iyi modelin belirlenmesini amaçlamaktadır. Projede hem dengeli hem dengesiz veri setleri üzerinde çalışılmış. Kullanıcı dostu bir arayüz ile model seçimi ve tahmin işlemleri kolaylaştırılmıştır.

Veri Setleri: Projede kullanılan veri setleri:

- **Orijinal Veri Seti:** diabetes.csv
- **Dengesiz Veri Seti:** diabetes_undersampled.csv

Veri seti 768 veri içermektedir. Veri setleri, Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction ve Age değişkenlerini içermektedir. Hedef değişken "Outcome" olup, 1 diyabet varlığı, 0 diyabet yokluğunu temsil etmektedir.

Orijinal veri setinde hedef değişken "Outcome" da 0 lar 500 ,1 ler 268 tanedir. Burada 1 lerin sayısını düşürüp 100 yaparak orijinal veri setini dengesiz veri seti yaptım. Bu veri setini “diabetes_undersampled.csv” kaydettim.

Kullanıcı Arayüzü Tasarımı: Ana ekran, modern ve kullanıcı dostu bir tasarıma sahiptir. Ekran bileşenleri:

- **Veri Yükleme:** Kullanıcı, orijinal ve dengesiz veri setlerinden birini yükleyebilir.
- **ComboBox (Model Seçimi):** KNN, Karar Ağacı ve Random Forest modellerinden birini seçebilir.
- **Model Seçimi (K- fold Butonu):** Orijinal veri seti üzerinde seçilen modelle göre k-fold uygulanması
- **Model Eğit Butonu:** Seçilen model eğitilir ve model başarımleri metrikleri gösterilir.
- **Tahmin Yap Butonu:** Eğitilen model ile tahmin yapılır ve predict_proba ve predict sonuçları ekranda gösterilir.

The screenshot shows the 'MainWindow' application interface. It features several input fields for features: Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction, and Age. There are also buttons for 'tahmin yap' (make prediction) and 'Smote'. A dropdown menu for 'Diyabet veri seti' (Diabetes data set) is set to '0'. Below this, there are buttons for 'Null %', 'Beri Doldurma' (Fill missing), 'Doldur' (Fill), 'Min-Max Norm', 'Normalizasyon' (Normalization), and 'KNN'. A 'K-fold' button is also present. A table displays the results of the model's predictions, showing the predicted outcome (0 or 1) and the predicted probability (predict_proba) for each instance.

	Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	DiabetesPedigreeFunction
1	6	148	72	35	0	33.6	0.627
2	1	85	66	29	0	26.6	0.351
3	8	183	64	0	0	23.3	0.672
4	1	89	66	23	94	28.1	0.167
5	0	137	40	35	168	43.1	2.288
6	5	116	74	0	0	25.6	0.201

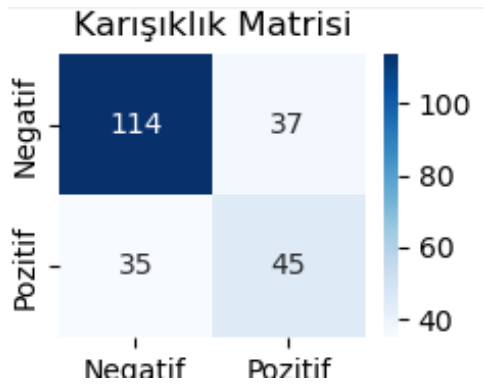
Kullanılan Modeller: Projede kullanılan modeller:

- KNN
- Decision Tree Classifier (Karar Ağacı)
- Random Forest Classifier
- K-Fold Çapraz Doğrulama

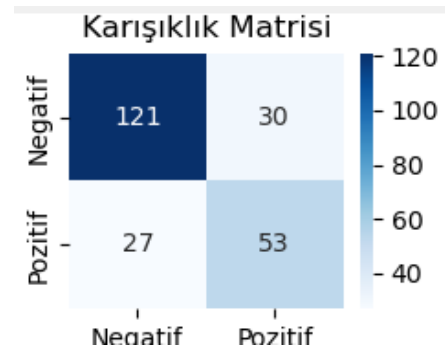
1)Orjinal Veri Seti Üzerinde Yapılan Çalışmalar

ComboBox dan seçile 3 metot ile model eğitimi gerçekleştirir. Bu eğitilen modelin doğruluk (Accuracy), duyarlılık (Recall), özgüllük (Specificity) ,F1 Skoru ve karmaşıklık matrisi sonuçları kullanıcıya verilmiştir.

1)Knn Uygulanması ve sonuçları :

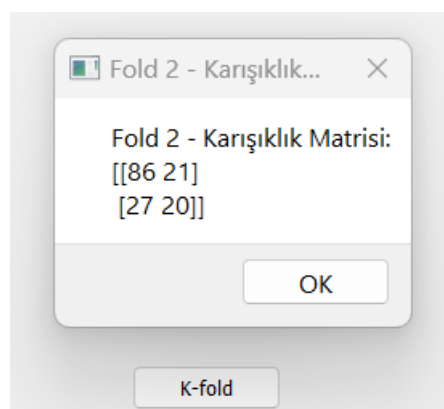
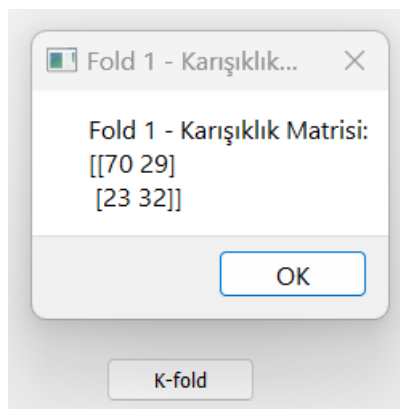


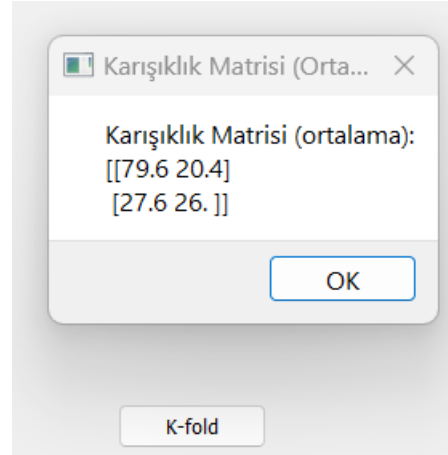
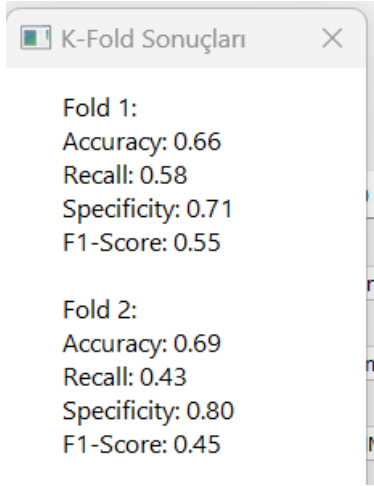
2)Random Forest Modeli ve Sonucu:



Orijinal veri seti üzerinde k-fold uygulanabilir. K-Fold Çapraz Doğrulama, veri setini k eşit parçaya böler ve her seferinde bir parçasını test veri seti olarak ayırırken kalan k-1 parçayı model eğitimi için kullanır.

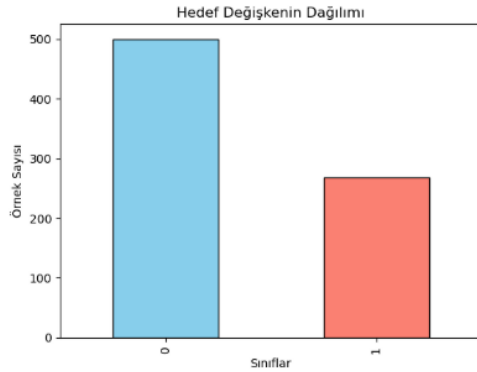
Her bir k (5) değeri için ayrı ayrı başarı metrikleri ve karmaşıklık matrisi oluşturur.



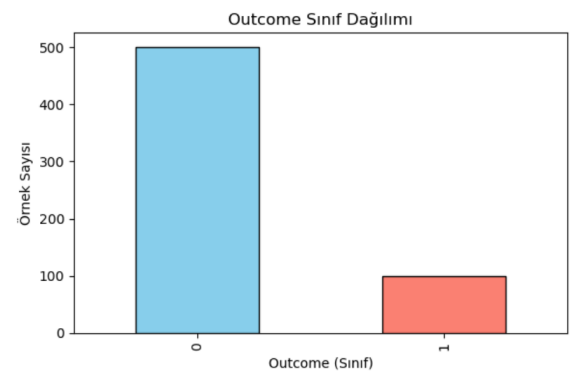


2) Dengesiz Veri seti üzerinde yapılan çalışmalar

Orijinal veri seti üzerinde yapılan işlem sonucu elde edilen dengesiz veri seti üzerinde ComboBox üzerinde seçilen 3 model ile çalışılabilir. Ve bu modellerin metrikleri ve karışıklık matrisi kullanıcıya verilir.



Şekil 1-Orijinal veri seti hedef değişken



Şekil 2-Dengesiz veri setindeki hedef değişken

3)Gürültülü Veri seti ile Çalışma

Gürültülü veriler, makine öğrenmesi modellerinin performansını olumsuz etkileyebilir ve yanlış tahminlere yol açabilir. Gürültüyü azaltmak için veri temizleme ve filtreleme teknikleri uygulanır.

Kullanıcı orijinal veri seti üzerinde istediği yüzde de null değer atayarak gürültülü veri seti elde eder. Bu gürültülü veri seti üzerinde comboBoxdan aldığı yöntemle göre ileri ile doldurma, geri ile doldurma ve ortalama ile doldurma seçenekleri ile gürültülü olan veri setini düzeltir. Düzeltelen veri seti üzerinde knn, random forest ve karar ağacı modellerinden biri seçilerek model eğitilir ve bunun sonuçları kullanıcıya verilir. (Başarı metrikleri ve karmaşıklık matrisi)

Diyabet veri seti

50

Null %

İleri Doldurma

Doldur

TextLabel

Min-Max Norm

Normalizasyon

TextLabel

KNN

Model Çalıştırma

TextLabel

Pregnancies	Glucose	BloodPressure	SkinThickness	Insulin	BMI	etesPedigre
		72.0		0.0		
			29.0	0.0	26.6	
8.0		64.0			23.3	0.672
1.0	89.0		23.0			
0.0	137.0	40.0				
5.0	116.0					

4)Normalizasyon Deneyi

Normalizasyonun amacı, veri setindeki değişkenlerin farklı ölçeklerde olmasının modellerin performansını olumsuz etkilemesini önlemektir. Normalizasyon, tüm değişkenleri aynı ölçek aralığına getirerek (genellikle [0,1] veya z-skor standardizasyonu ile) modelin daha dengeli ve doğru tahminler yapmasını sağlar.

Orijinal veri setindeki hedef değişken “Outcome” sınıfı dışında kalan değişkenler üzerinde normalizasyon uygulanır. İki çeşit normalizasyon uygulanabilir. Min-max normalizasyonu ve z-skor standardizasyonu uygulanabilir.

Kullanıcı seçtiği üç modelden biri ile eğitim gerçekleştirebilir ve bunun başarı metrikleri ve karmaşıklık matrisi kullanıcıya verilir.

Diyabet veri seti

0

Null %

İleri Doldurma

Doldur

TextLabel

Min-Max Norm

Normalizasyon

Accuracy: 0.75Precision: 0.64Recall: 0.66F1-S

KNN

Model Çalıştırma

TextLabel

	Glucose	BloodPressure	SkinThickness	Insulin	BN
1	0.7437185929648241	0.5901639344262295	0....	0.0	0....
2	0.4271356783919598	0.5409836065573771	0....	0.0	0....
3	0.9195979899497487	0.5245901639344263	0.0	0.0	0....
4	0.4472361809045226	0.5409836065573771	0....	0....	0....
5	0.6884422110552764	0....	0....	0....	0....
6	0.5820145728662216	0.6065572770401802	0.0	0.0	0

5)Kullanıcının Girdiği değerlere Göre Hedef sınıfı Tahmin Etme

Pregnancies, Glucose, BloodPressure, SkinThickness, Insulin, BMI, DiabetesPedigreeFunction ve Age değişkenlerini kullanıcıdan alınır ve Outcome sınıfı tahmin edilir.

Pregnancies	6
Glucose	148
BloodPressure	72
SkinThickness	35
Insulin	0
BMI	33.6
DiabetesPedigreeFunction	0.627
Age	50
<button>tahmin yap</button>	
Sonuç: 1 Olasılık (0): 0.13 Olasılık (1): 0.87	

Pregnancies	1
Glucose	85
BloodPressure	66
SkinThickness	29
Insulin	0
BMI	26.6
DiabetesPedigreeFunction	0.351
Age	31
<button>tahmin yap</button>	
Sonuç: 0 Olasılık (0): 0.94 Olasılık (1): 0.06	

Tahmin (predict metodu): Model, kullanıcıdan alınan özellikler ile predict metodunu kullanarak hedef değişkenin tahminini gerçekleştirdi. Örneğin, kullanıcının girdiği verilere göre model hedef değişkenin 0 mı yoksa 1 mi olacağını tahmin etti.

Olasılık (predict_proba metodu): predict_proba metodu ise her iki sınıf için (0 ve 1) olasılıkları hesapladı. Bu, modelin tahminin doğruluğu hakkında daha fazla bilgi sağlar; örneğin, modelin tahmin ettiği 1 sınıfının olasılığı %75, 0 sınıfının olasılığı ise %25 olabilir. Bu olasılıklar, modelin karar verme sürecine dair güven seviyesini yansıtır.