

name	MDVP:Fo(Hz)	MDVP:Fhi(Hz)	MDVP:Flo(Hz)	MDVP:Jitter(%)	MDVP:Jitter(Abs)	MDVP:RAP	MDVP:PPQ
phon_R01_S01_1	119.992	157.302	74.997	0.00784	0.00007	0.00370	0.00554
phon_R01_S01_2	122.400	148.650	113.819	0.00968	0.00008	0.00465	0.00696
phon_R01_S01_3	116.682	131.111	111.555	0.01050	0.00009	0.00544	0.00781
phon_R01_S01_4	116.676	137.871	111.366	0.00997	0.00009	0.00502	0.00698
phon_R01_S01_5	116.014	141.781	110.655	0.01284	0.00011	0.00655	0.00908

count 195.000000

mean 0.753846

std 0.431878

min 0.000000

25% 1.000000

50% 1.000000

75% 1.000000

max 1.000000

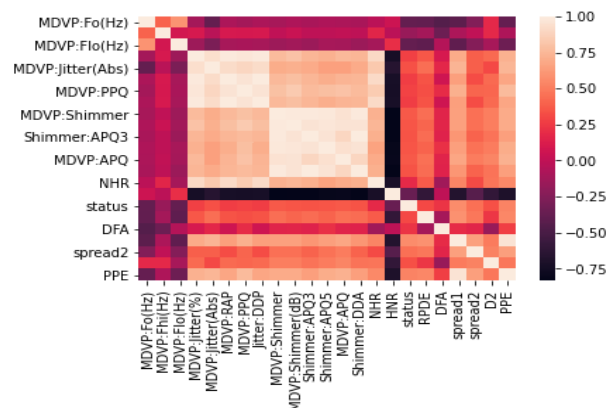
Name: status, dtype: float64

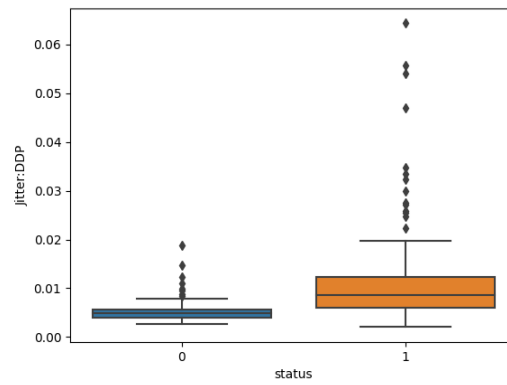
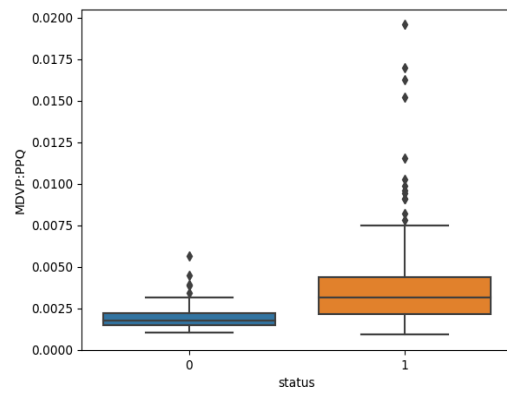
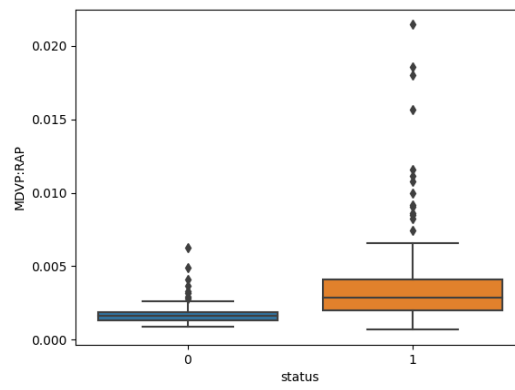
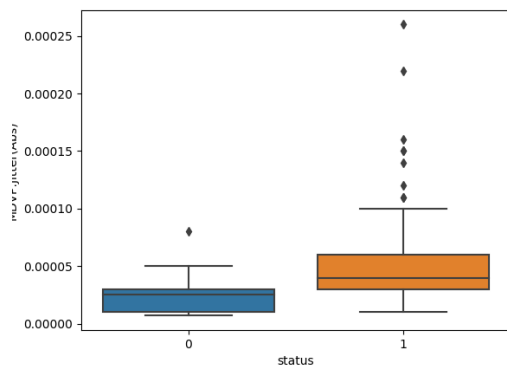
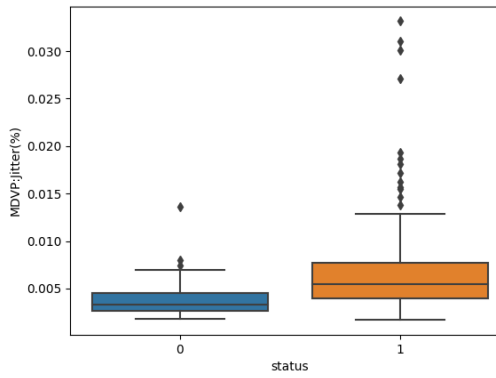
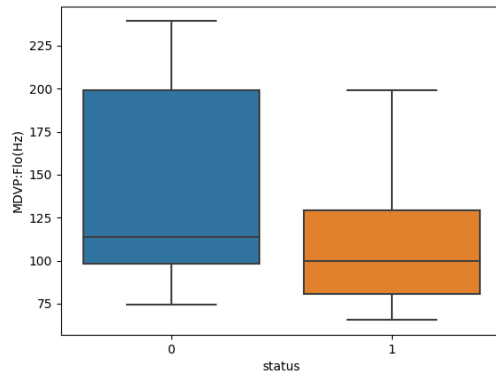
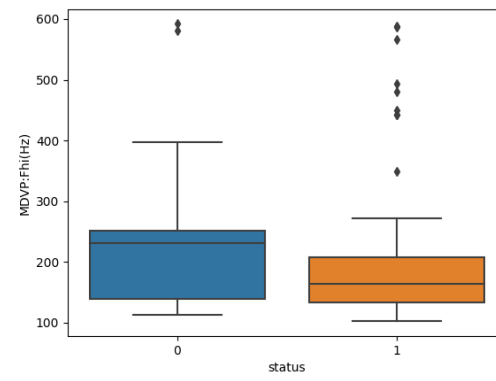
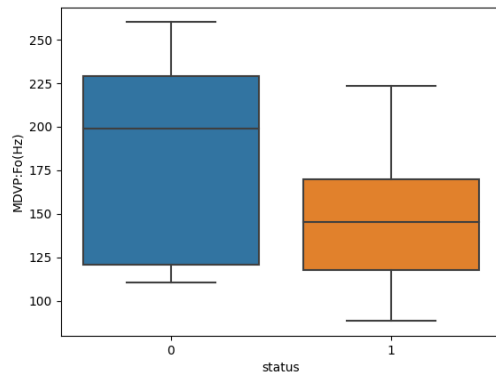
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 195 entries, 0 to 194
Data columns (total 24 columns):
#   Column                Non-Null Count  Dtype
---  -
0   name                  195 non-null    object
1   MDVP:Fo(Hz)          195 non-null    float64
2   MDVP:Fhi(Hz)         195 non-null    float64
3   MDVP:Flo(Hz)         195 non-null    float64
4   MDVP:Jitter(%)       195 non-null    float64
5   MDVP:Jitter(Abs)     195 non-null    float64
6   MDVP:RAP             195 non-null    float64
7   MDVP:PPQ             195 non-null    float64
8   Jitter:DDP           195 non-null    float64
9   MDVP:Shimmer         195 non-null    float64
10  MDVP:Shimmer(dB)     195 non-null    float64
11  Shimmer:APQ3         195 non-null    float64
12  Shimmer:APQ5         195 non-null    float64
13  MDVP:APQ             195 non-null    float64
14  Shimmer:DDA          195 non-null    float64
15  NHR                  195 non-null    float64
16  HNR                  195 non-null    float64
17  status               195 non-null    int64
18  RPDE                 195 non-null    float64
19  DFA                  195 non-null    float64
20  spread1              195 non-null    float64
21  spread2              195 non-null    float64
22  D2                   195 non-null    float64
23  PPE                  195 non-null    float64
dtypes: float64(22), int64(1), object(1)
memory usage: 36.7+ KB
```

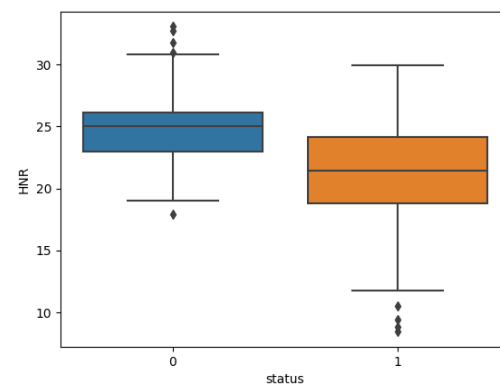
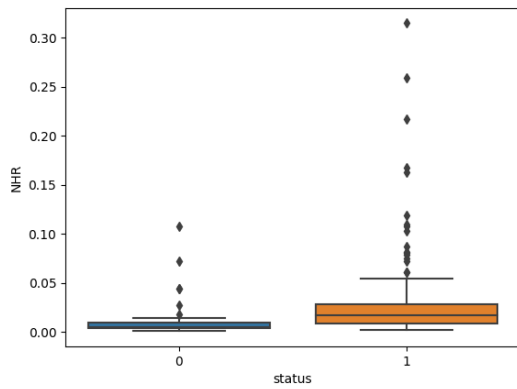
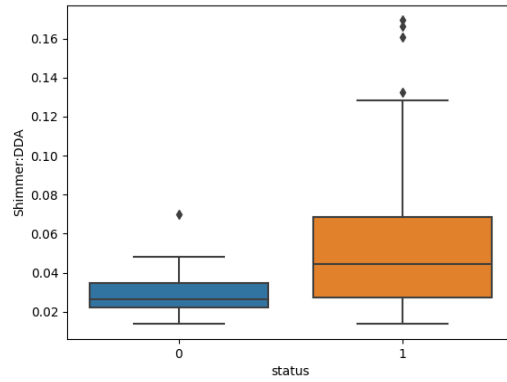
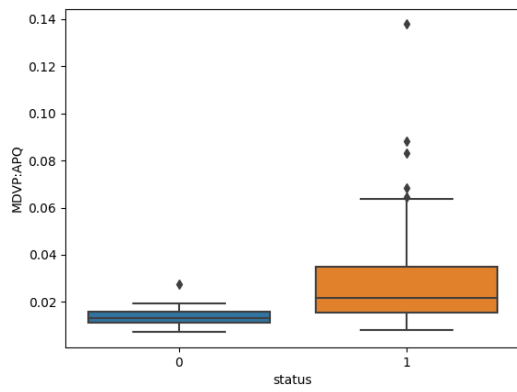
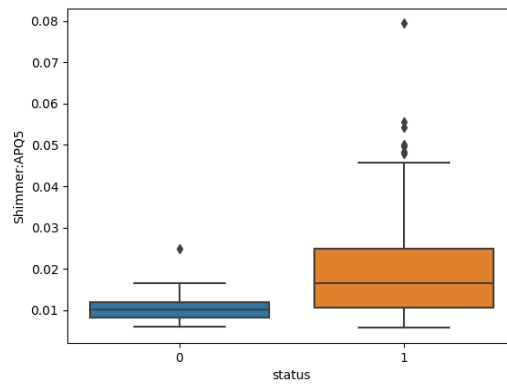
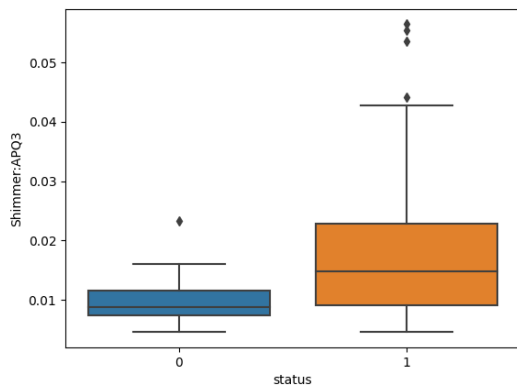
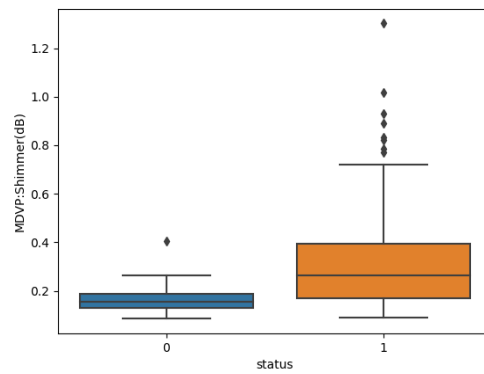
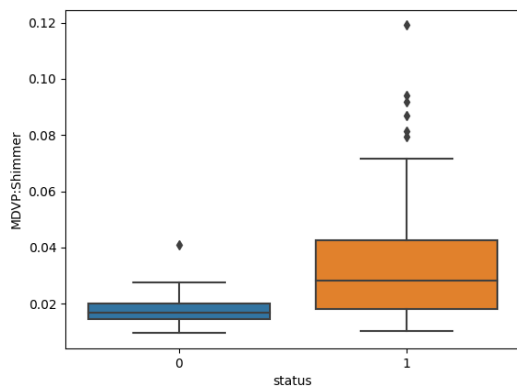
Bu noktada verilerimize genel olarak baktığımızda oldukça sağlıklı bir veri seti elde etmiş olduğumuzu görüyoruz. Herhangi bir outlier durumu görünmemekle birlikte genel olarak hastaların yoğun olduğu bir veri seti karşımızda bulunmakta. Bu aslında algoritmayı sağlıklı ve hastalıklı olarak ayırmada zorlayacak bir durum. Eşit dağılım olsaydı, algoritmalarımız daha kesin çalışırdı fakat yine de çok iyi sonuçlar aldığımızı belirtmem gerek. 195 adet veriden fena olmayan sonuçlar aldık fakat verilerimizin sayısı daha çok olsaydı tahminimce daha iyi eğitilmiş bir algoritmalarımız olabilirdi.

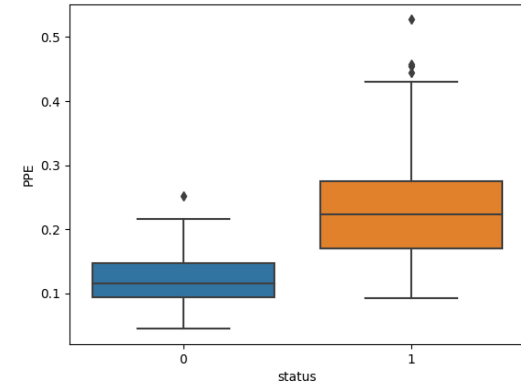
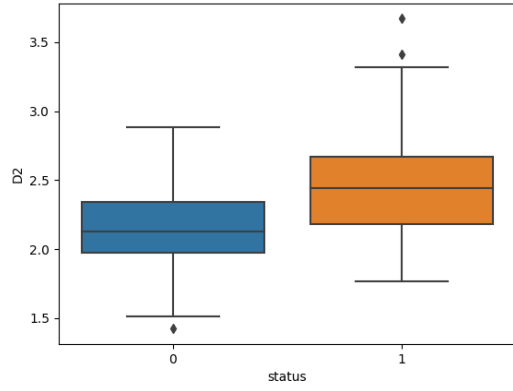
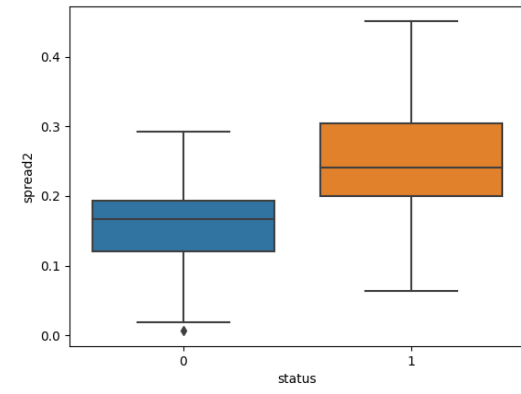
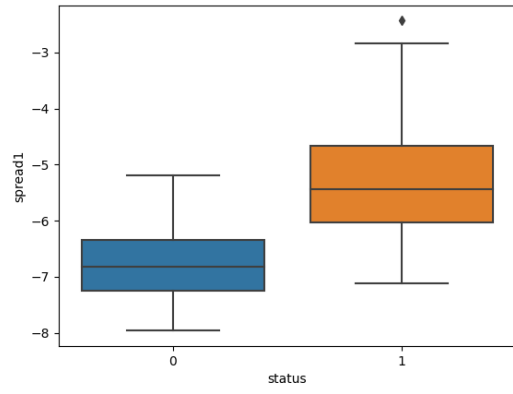
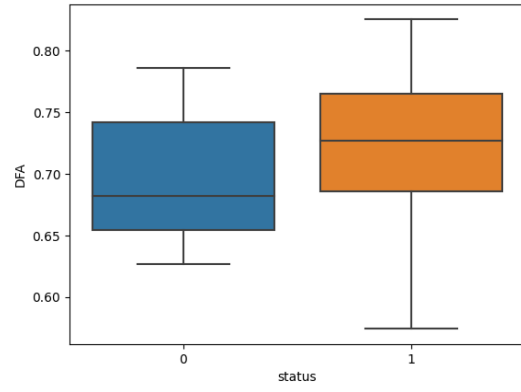
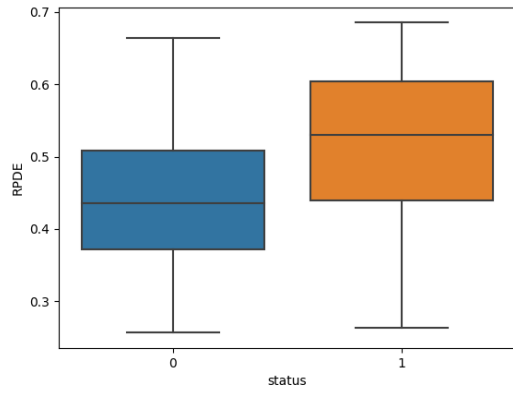
Verilerin görsel ilişkilerine bir göz gezdirelim

```
sns.heatmap(df.corr());
# data üzerinde korrelasyon spread1,spread2 ve ppe featureslerinde yoğunlukta heatmap üzerinden görebiliriz.
```

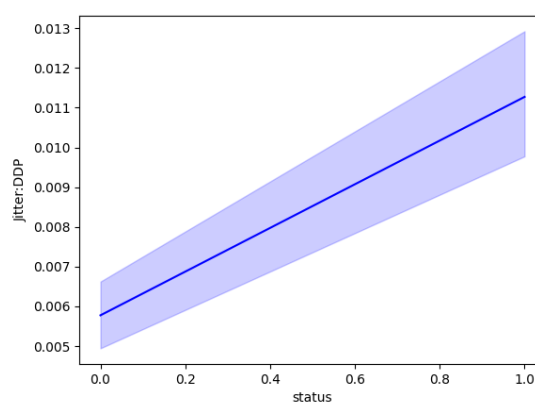
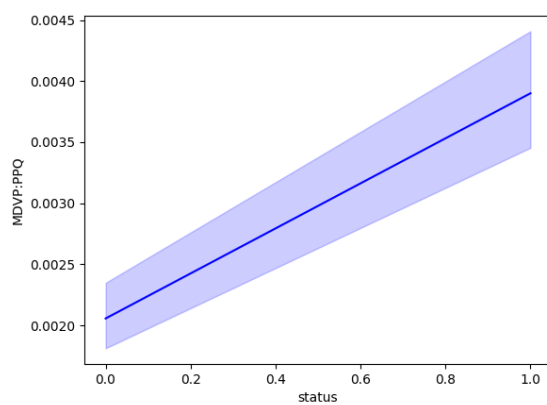
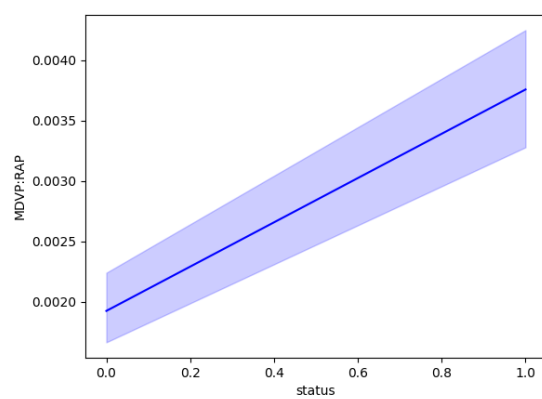
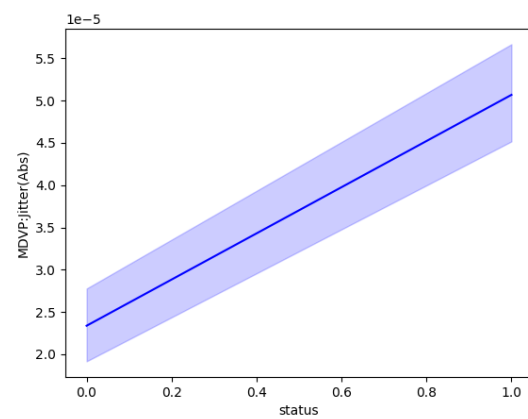
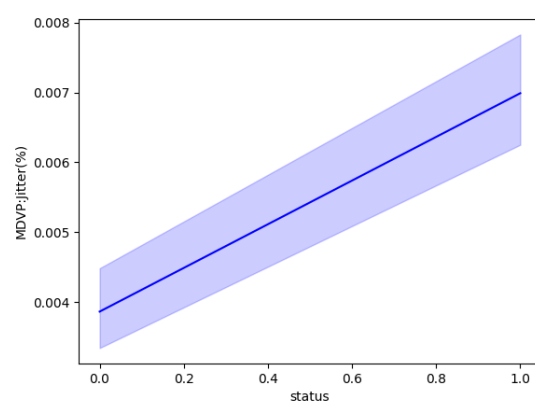
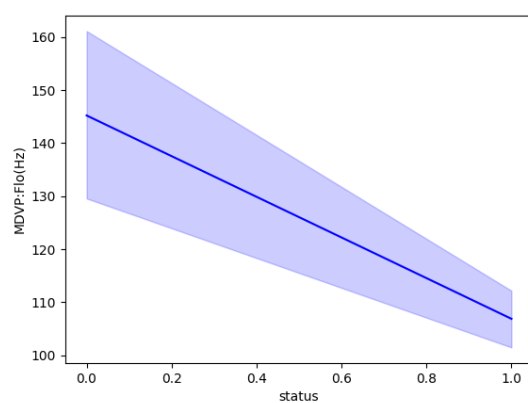
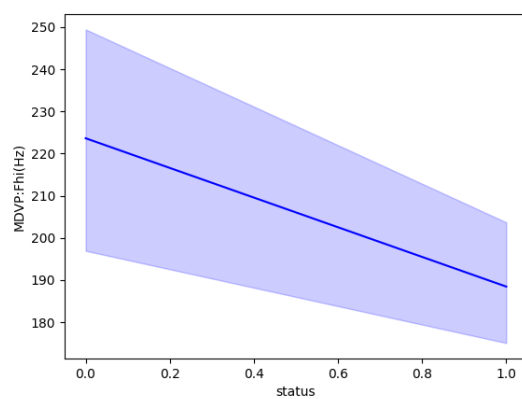
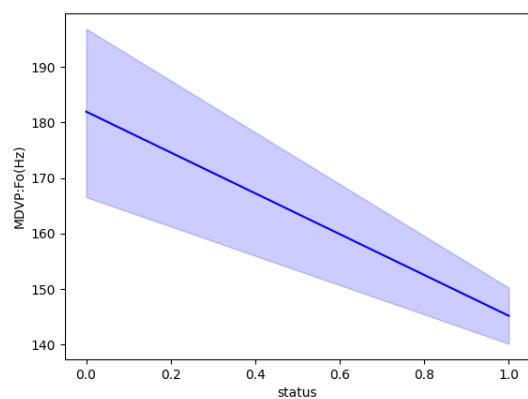


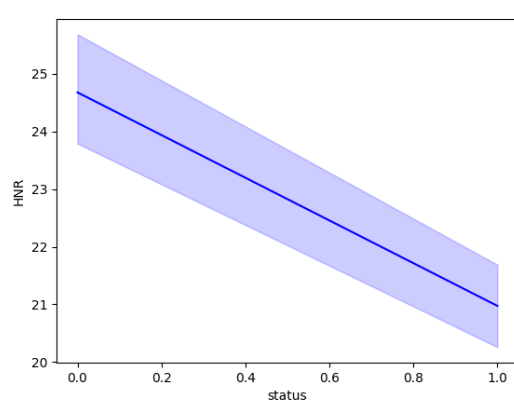
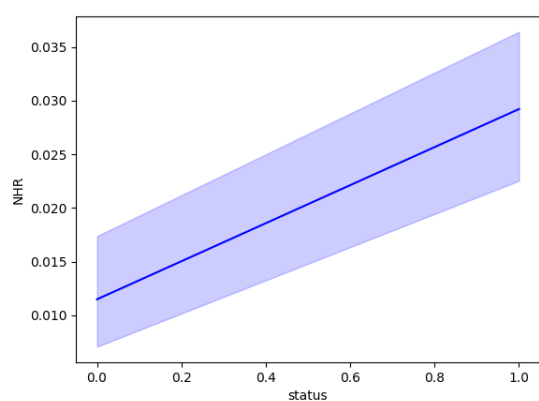
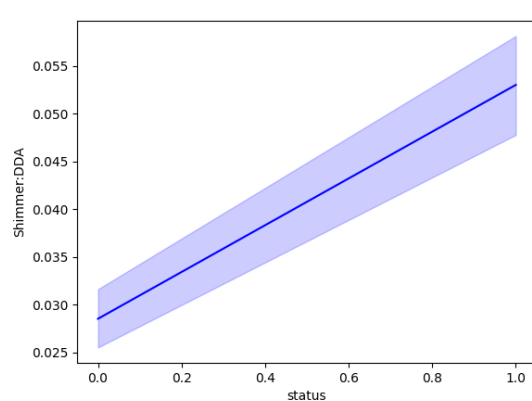
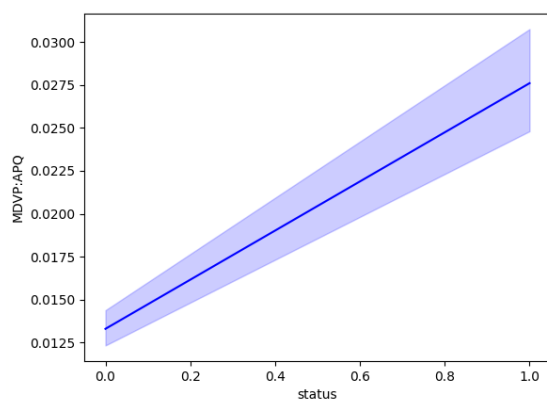
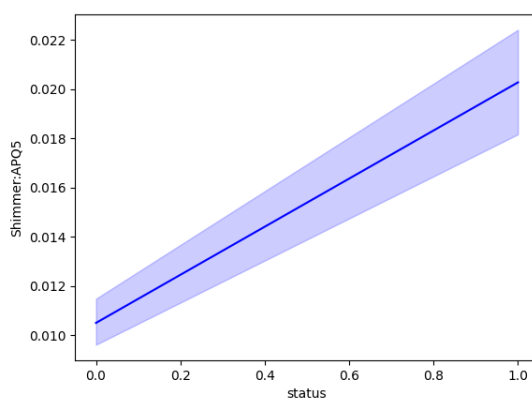
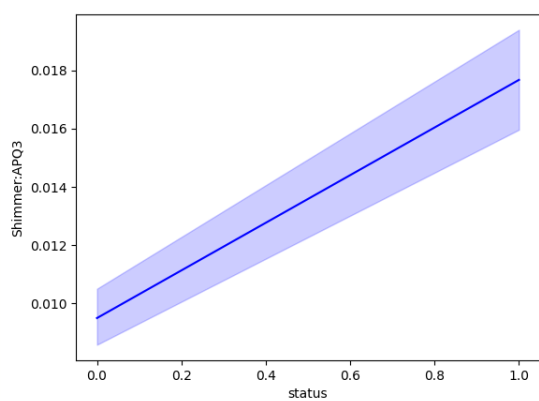
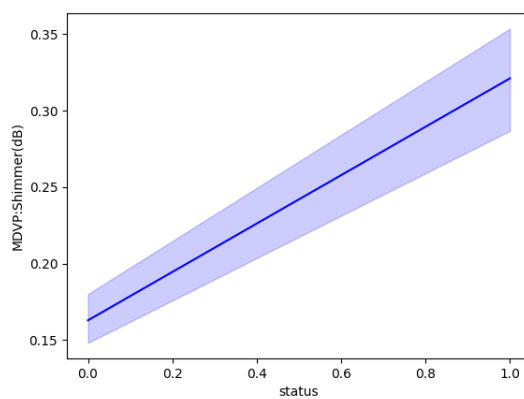
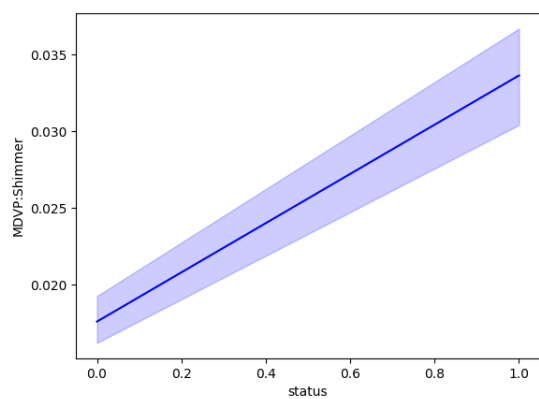


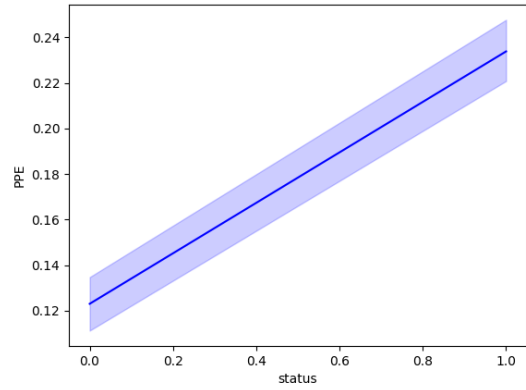
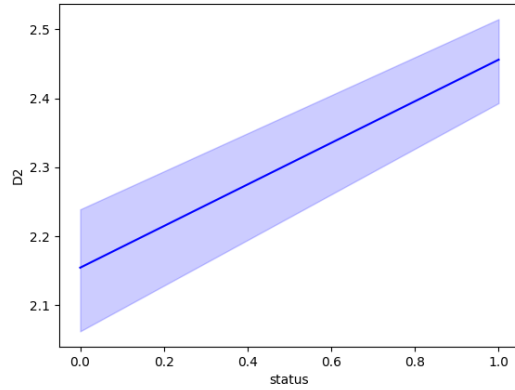
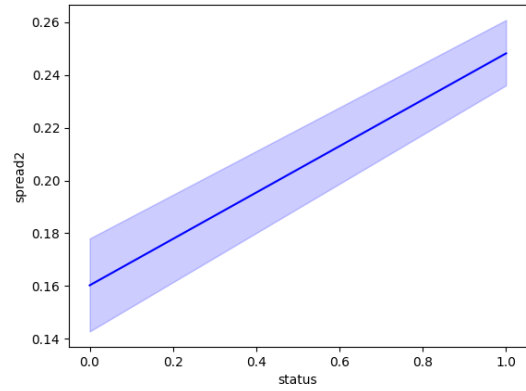
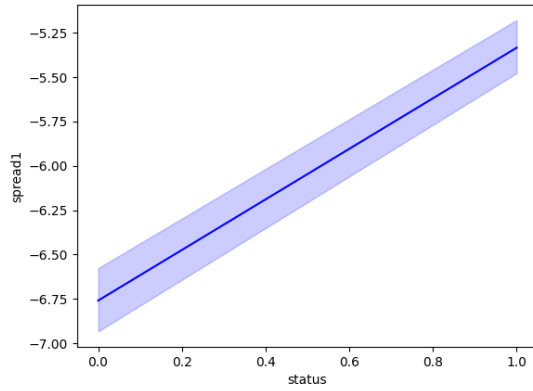
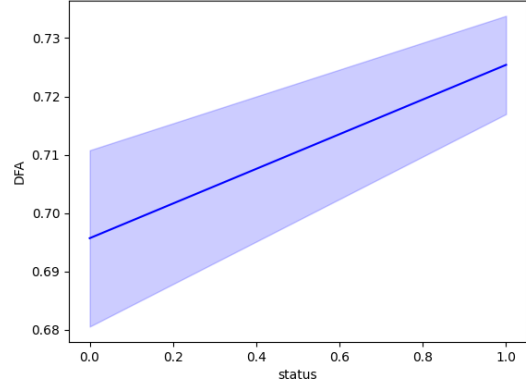
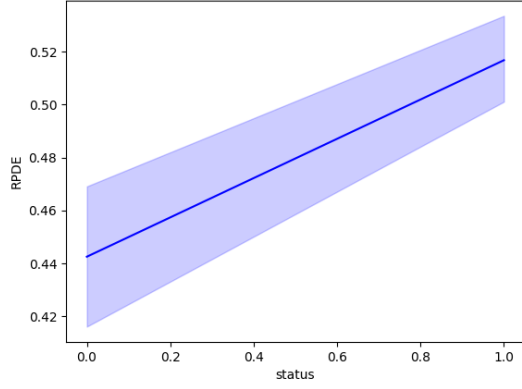




spread1, spread2 ve ppe featureleri güzel ayırım yapmış. Değerleri yüksek olanlar genelde hasta, düşük olanlar ise genelde hasta değil çıkarmını yapabiliriz bu noktada.







MDVP-FO(Hz),MDVP-Fhi(Hz), MDVP-Flo(Hz) ve HNR değerlerinin düşmesi demek hastalıklı olmanız demektir. Bu değerlerin yüksek olması hastalığı ayıran en net etkenlerden biri olarak görülmektedir.

Verilerimizi Standart şekilde scale ettikten sonra preprocess işlemlerimizi bitirmiş olacağız. Bundan sonra gidişata göre train,test ayrımı yapacağımızı söylemiştim. Deneme yanılma yoluyla uygun ayrımı 75 train, 25 test şeklinde karar verdim. Seçtiğimiz 4 modeli sözlük veri yapısında saklayıp ilgili eğitimini gerçekleştirdik. Verilerimiz hazır olduğuna göre karşılaştırma kısmına geçebiliriz.

XGBoost accuraty is: % 95.91836734693877

Logistic_Regression accuraty is: % 81.63265306122449

SVM_Classifier accuraty is: % 89.79591836734694

Desicion_Tree accuraty is: % 83.6734693877551

- XGBoost 100 deneme sonucu: 93-97 arasında bir sonuç verdi. Genel olarak kendisi mükemmel bir algoritma olduğu için az veriye rağmen overfitting yapmamaya gayret ediyor. Kendisi maksimum shuffle edilen test verilerinde 3 kere hata yaptı.

- Logistic Regression 100 deneme sonucu: Bazı noktalarda accuraty, %90'ların üstüne çıktı fakat genel olarak 80-84 arasında geziyor. Bununla birlikte yüksek accuraty değerleri ile düşük accuraty değerleri arasındaki fark, düşük veriden kaynaklı overfitting göstergesi olarak yorumluyoruz.

- SVM Classifier 100 deneme sonucu: Bu noktada SVM makinesi, maksimum %91 minimum %84 değerlerinde dolaştı. Burada lojistik regresyondan daha iyi bir sonuç ama genel olarak overfitting durumu görülmekte, eğer C değerini overfitting için yükseltirsek accuraty seviyesi %70'lere düşüp az verinin cezası olarak tahminde fazla hata yapma durumunda kalıyor. Bu noktada varyans arttırılması, daha genel bir öğrenme yaparken az veri için sınıflandırma hatasına neden oluyor. Yani az veri konusunda hala XGBoost 1 numara.

- Desicion Tree 100 deneme sonucu: Karar ağaçları, bazen XGBoost accuraty değerleri yakalasa da bazen de lojistik regresyon kadar düşük değerler verebiliyor. Bu noktada maximum 93, minimum 81 accuarty aralığı seçebiliriz. Karar ağaçlarınının değişken yapısı, random-forest tekniğinden dolayı bazen iyi bir çizgi çekerken bazen overfitting durumunda kalıyor. Bu noktada dengesiz bir öğrenim söz konusu diyebilirim. Bunun sebebi ise bariz, verilerimiz az sayıda.

XGBoost Model Classification predict -- Real Data is Logistic Regression Model Classification predict -- Real Data is

[illegible]

Total Mistake is : 2

[illegible]

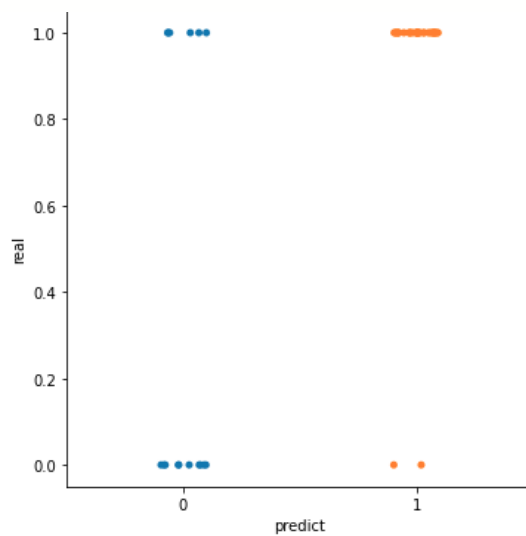
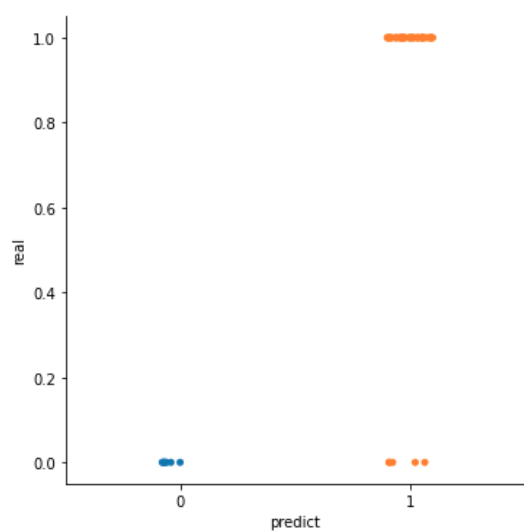
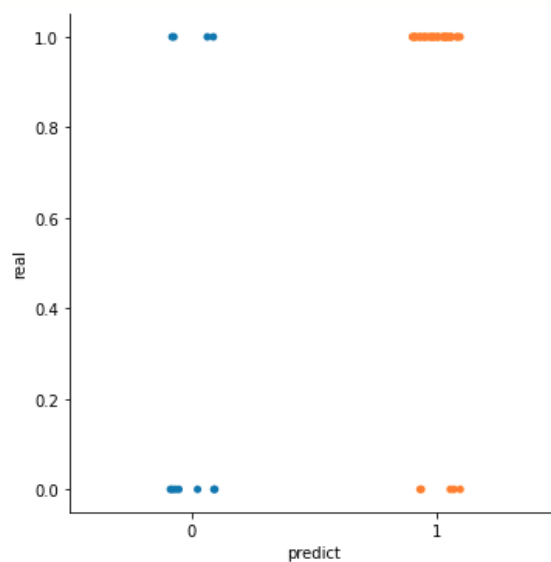
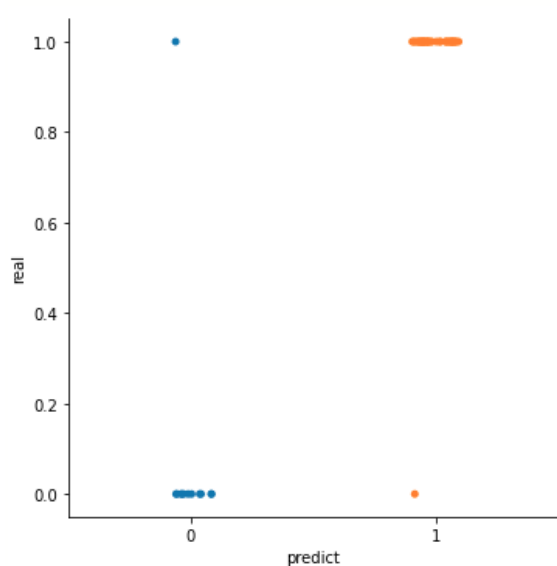
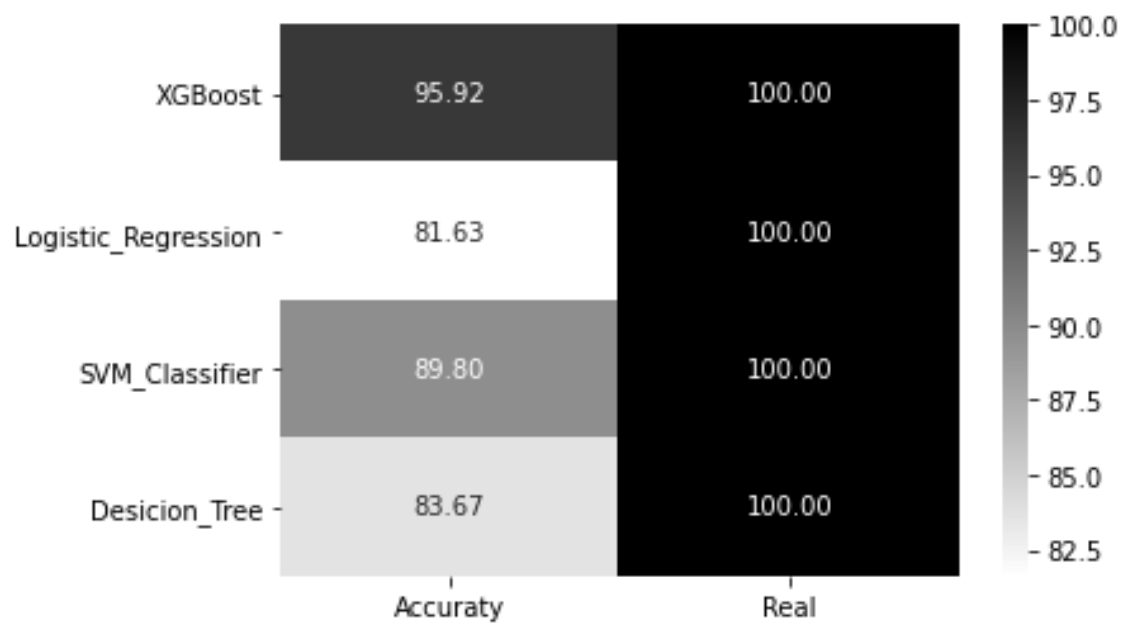
Total Mistake is : 9

SVM_Classifier Model Classification predict -- Real Data is Desicion_Tree Model Classification predict -- Real Data is

[illegible]

Total Mistake is : 5

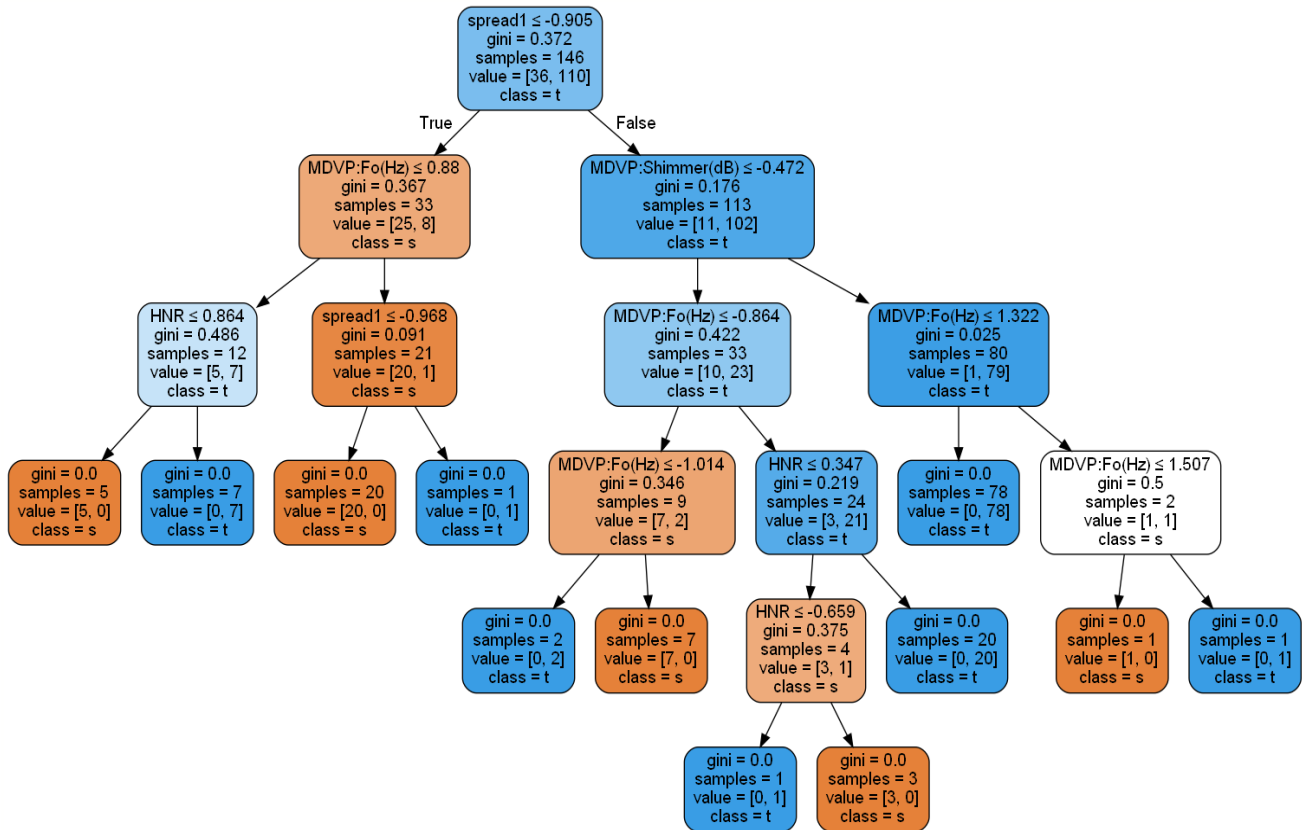
Total Mistake is : 8



- Sonuç: Düşük veri arasında modeller için karşılaştırma yaptığımızda bir sıralama yaparsak;

- 1-) XGBoost
- 2-) SVM Machine
- 3-) Desicion Tree
- 4-) Logistic Regression

Sıralamasını elde etmiş bulunmaktayız.



Desicion Tree için örnek bir tree oluşumunu da bu şekilde görmekteyiz. Desicion Tree için gini yerine cross-entropy seçtiğimiz zaman genel accuraty düşüklüğü elde etmiş bulunmaktayız. Az verilerimizden dolayı cross-entropy, çok keskin kararlar veriyor ve birbirlerine yakın hasta-sağlıklı dataları için sınıflandırma hatası yapıyor. Bu yüzden gini, cost fonksiyonu olarak uygun görülmektedir. Bu projede, hyperparameter optimization konusu dâhil değildir. Bu yüzden parametreler, deneme yanılma yoluyla elde edilmiştir.