

Parkinson's Disease prediction from Voice signals

Created By: Mrinal Bhan & Gautam Gupta

Parkinson's disease

- ❑ Abnormalities of the Parkinson's disease speech can be associated with several dimensions
- ❑ Symptoms include impairment in the normal production of vocal sounds (dysphonia), and problems with normal articulation of speech (dysarthria)
- ❑ Dysphonic symptoms typically include
 - ❑ reduced loudness,
 - ❑ breathiness,
 - ❑ roughness,
 - ❑ decreased energy in the higher parts of the harmonic spectrum, and
 - ❑ exaggerated vocal tremor

Speech measurement for PD voice disorder:

- ❑ Traditional methods:
 - ❑ Fundamental frequency
 - ❑ Absolute sound pressure level
 - ❑ Jitter
 - ❑ Shimmer
 - ❑ Noise-to-harmonics ratios
- ❑ Novel measurement methods (based on non-linear dynamic systems):
 - ❑ Recurrence period density entropy (RPDE), D₂
 - ❑ Detrended fluctuation analysis (DFA)
 - ❑ Pitch period entropy (PPE)
 - ❑ Spread 1, Spread 2 – Nonlinear measures of fundamental frequency variation

References:

1. Little MA, McSharry PE, Hunter EJ, Spielman J, Ramig LO. Suitability of dysphonia measurements for telemonitoring of Parkinson's disease. IEEE Trans Biomed Eng. 2009 Apr;56(4):1015. doi: 10.1109/TBME.2008.2005954. PMID: 21399744; PMCID: PMC3051371.
2. Rusz J, Cmejla R, Ruzickova H, Ruzicka E. Quantitative acoustic measurements for characterization of speech and voice disorders in early untreated Parkinson's disease. J Acoust Soc Am. 2011 Jan;129(1):350-67. doi: 10.1121/1.3514381. PMID: 21303016.

Objective:

- ❑ Find the key predictors, especially speech-related characteristics, of PD
- ❑ Try at least three different machine learning approaches to PD identification and Find the best approach.

EDA

- ❑ We have 195 rows and 24 variables.
- ❑ No null values or missing values in the data
- ❑ All of the independent variables are numeric

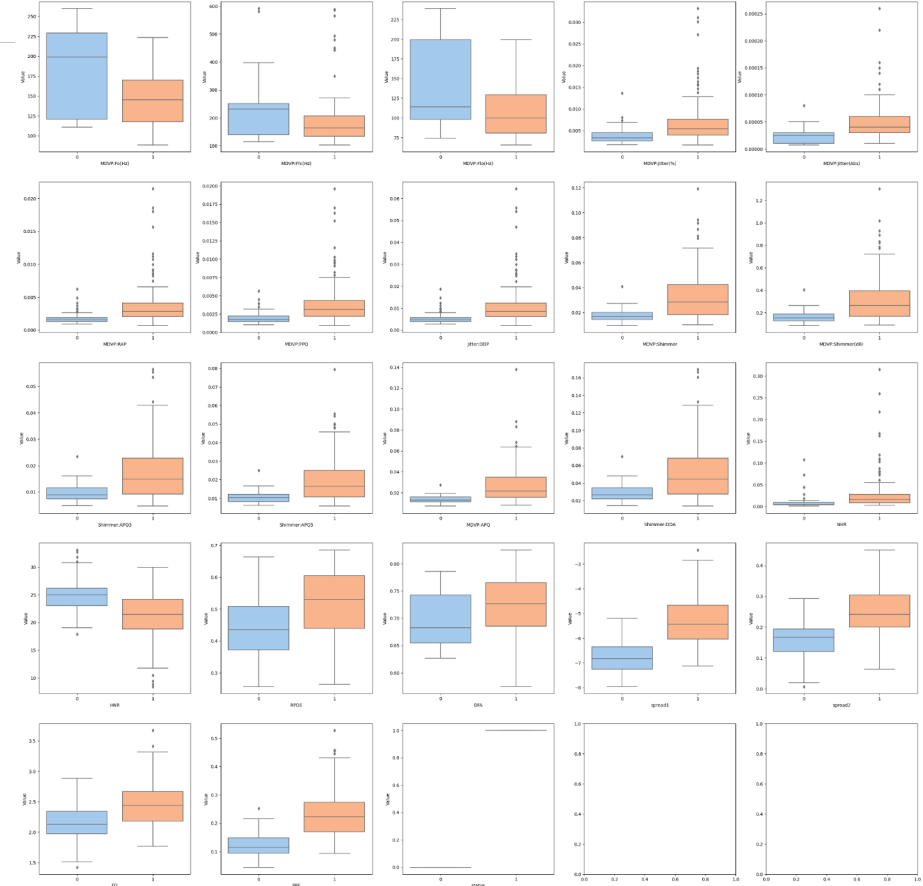
`data.shape`

`(195, 24)`

<code>data.isna().sum().sort_values</code>		<code>data.dtypes</code>	
<code>name</code>	<code>0</code>	<code>name</code>	<code>object</code>
<code>MDVP:Fo(Hz)</code>	<code>0</code>	<code>MDVP:Fo(Hz)</code>	<code>float64</code>
<code>PPE</code>	<code>0</code>	<code>MDVP:Fhi(Hz)</code>	<code>float64</code>
<code>D2</code>	<code>0</code>	<code>MDVP:Flo(Hz)</code>	<code>float64</code>
<code>spread2</code>	<code>0</code>	<code>MDVP:Jitter(%)</code>	<code>float64</code>
<code>spread1</code>	<code>0</code>	<code>MDVP:Jitter(Abs)</code>	<code>float64</code>
<code>DFA</code>	<code>0</code>	<code>MDVP:RAP</code>	<code>float64</code>
<code>RPDE</code>	<code>0</code>	<code>MDVP:PPQ</code>	<code>float64</code>
<code>HNR</code>	<code>0</code>	<code>Jitter:DDP</code>	<code>float64</code>
<code>NHR</code>	<code>0</code>	<code>MDVP:Shimmer</code>	<code>float64</code>
<code>Shimmer:DDA</code>	<code>0</code>	<code>MDVP:Shimmer(dB)</code>	<code>float64</code>
<code>MDVP:APQ</code>	<code>0</code>	<code>Shimmer:APQ3</code>	<code>float64</code>
<code>Shimmer:APQ5</code>	<code>0</code>	<code>Shimmer:APQ5</code>	<code>float64</code>
<code>Shimmer:APQ3</code>	<code>0</code>	<code>MDVP:APQ</code>	<code>float64</code>
<code>MDVP:Shimmer(dB)</code>	<code>0</code>	<code>Shimmer:DDA</code>	<code>float64</code>
<code>MDVP:Shimmer</code>	<code>0</code>	<code>NHR</code>	<code>float64</code>
<code>Jitter:DDP</code>	<code>0</code>	<code>HNR</code>	<code>float64</code>
<code>MDVP:PPQ</code>	<code>0</code>	<code>RPDE</code>	<code>float64</code>
<code>MDVP:RAP</code>	<code>0</code>	<code>DFA</code>	<code>float64</code>
<code>MDVP:Jitter(Abs)</code>	<code>0</code>	<code>spread1</code>	<code>float64</code>
<code>MDVP:Jitter(%)</code>	<code>0</code>	<code>spread2</code>	<code>float64</code>
<code>MDVP:Flo(Hz)</code>	<code>0</code>	<code>D2</code>	<code>float64</code>
<code>MDVP:Fhi(Hz)</code>	<code>0</code>	<code>PPE</code>	<code>float64</code>
<code>status</code>	<code>0</code>	<code>status</code>	<code>int64</code>
<code>dtype: int64</code>		<code>dtype: object</code>	

EDA

Number of outliers in MDVP:Fo(Hz): 0
 Number of outliers in MDVP:Fhi(Hz): 11
 Number of outliers in MDVP:Flo(Hz): 9
 Number of outliers in MDVP:Jitter(%): 14
 Number of outliers in MDVP:Jitter(Abs): 7
 Number of outliers in MDVP:RAP: 14
 Number of outliers in MDVP:PPQ: 15
 Number of outliers in Jitter:DDP: 14
 Number of outliers in MDVP:Shimmer: 8
 Number of outliers in MDVP:Shimmer(dB): 10
 Number of outliers in Shimmer:APQ3: 6
 Number of outliers in Shimmer:APQ5: 13
 Number of outliers in MDVP:APQ: 12
 Number of outliers in Shimmer:DDA: 6
 Number of outliers in NHR: 19
 Number of outliers in HNR: 3
 Number of outliers in RPDE: 0
 Number of outliers in DFA: 0
 Number of outliers in spread1: 4
 Number of outliers in spread2: 2
 Number of outliers in D2: 1
 Number of outliers in PPE: 5



- ☐ We see that there are outliers in many features.
- ☐ But all of them are in the possible biological ranges for those variables.
- ☐ Decided to let the outliers be as is in the data,

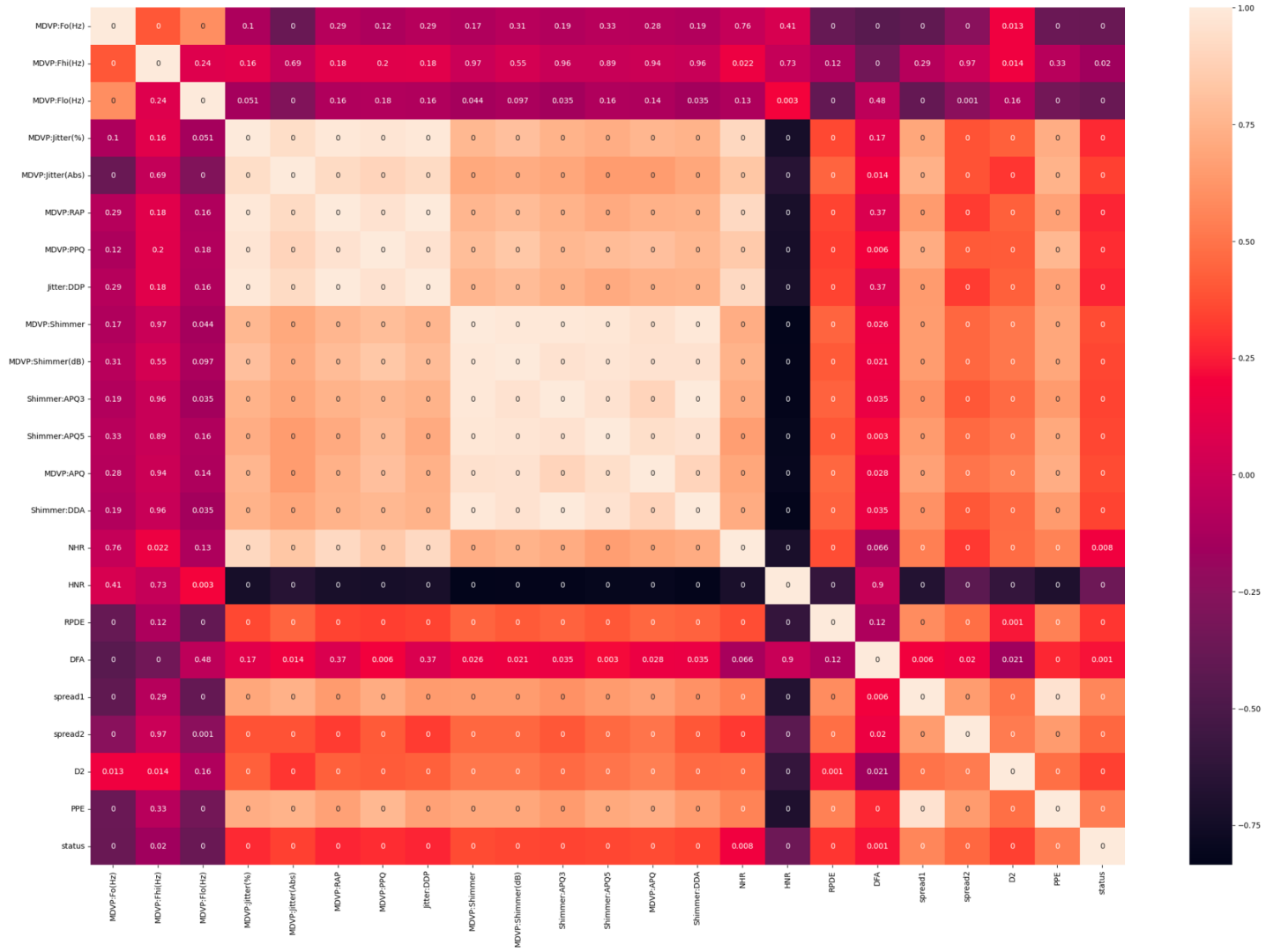
Correlation

❑ We can see clusters of strong correlation between the variables

❑ For ex, Measurements of Jitter are strongly correlated with each other and a decent correlation is seen between Shimmer and Jitter variables.

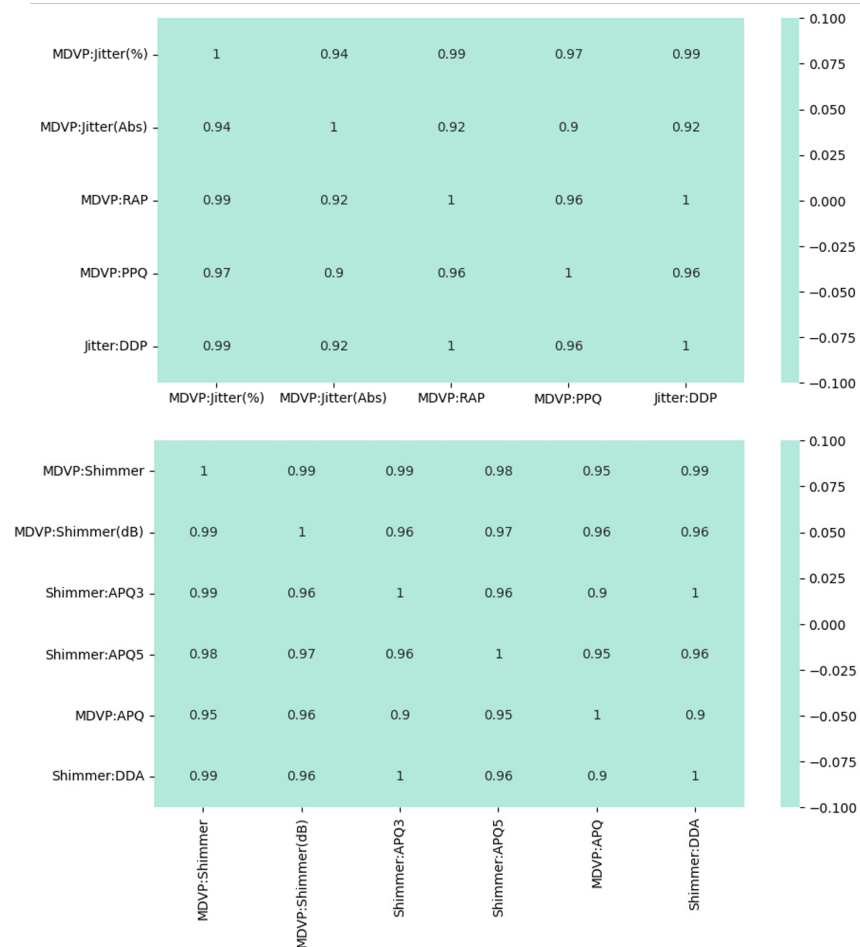
❑ Also, there is a strong negative correlation HNR and the rest of the variables.

❑ The novel measurements like RPDE, DFA, PPE etc are not correlated strongly as they are non-linear measures.



Checking correlation within clusters:

- ❑ As we saw earlier, There is a strong correlation between the different measures of Jitter and similarly in the Shimmer variables.
- ❑ The information from using all the variables is limited.
- ❑ Therefore, we can drop the variables or combine them using Principal component analysis.
- ❑ We tried both approaches



Variables:

- ❑ We decided to consider the range of the Fundamental frequency instead of the Max and Min values of the Fundamental frequency based on our research.

```
data['MDVP:F_Spd'] = data['MDVP:Fhi(Hz)'] - data['MDVP:Flo(Hz)']
```

- ❑ Rest of the variables were left as is
- ❑ There 12 independent variables left in the data we can use for modelling
- ❑ The remaining independent variables in the data are:

```
df_scaled.columns
```

```
Index(['MDVP:Fo(Hz)', 'MDVP:Jitter(Abs)', 'MDVP:Shimmer', 'NHR', 'HNR', 'RPDE',  
      'DFA', 'spread1', 'spread2', 'D2', 'PPE', 'MDVP:F_Spd'],  
      dtype='object')
```

- ❑ First, we went with dropping the correlated variables. We also did PCA in coming up slides

Models (1):

- ❑ Normalised the data using the standard scaler

- ❑ We have built 4 different models
 - a. Decision Tree
 - b. Random Forest
 - c. SVM
 - d. k-Neural network

- ❑ The performance of the models on the test data:

	Metric	DT	RF	SVM	KNN
0	Accuracy	0.820513	0.948718	0.846154	0.897436
1	F1-Score	0.881356	0.962963	0.892857	0.928571
2	Recall	0.962963	0.962963	0.925926	0.962963
3	Precision	0.812500	0.962963	0.862069	0.896552
4	R2-Score	0.157407	0.759259	0.277778	0.518519

Principal Component Analysis:

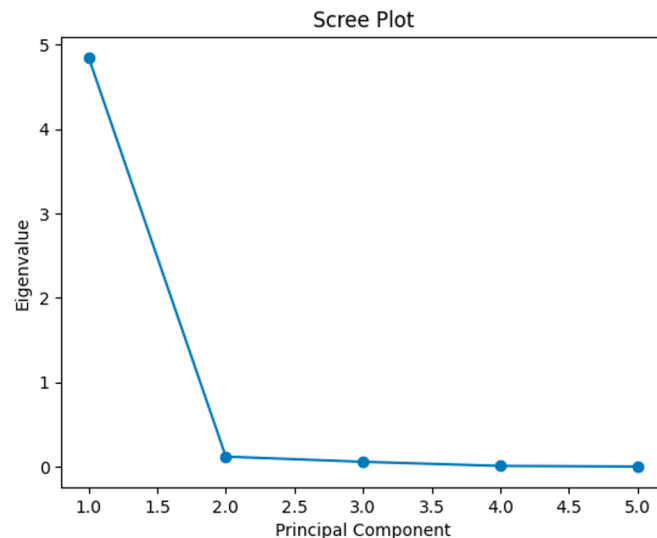
- ❑ As mentioned earlier, we have also applied PCA instead of dropping the correlated variables.
- ❑ We applied PCA on two different subsets of variables

```
variation_freq = ['MDVP:Jitter(%)', 'MDVP:Jitter(Abs)', 'MDVP:RAP', 'MDVP:PPQ', 'Jitter:DDP']  
variation_amp = ['MDVP:Shimmer', 'MDVP:Shimmer(dB)', 'Shimmer:APQ3', 'Shimmer:APQ5', 'MDVP:APQ', 'Shimmer:DDA']
```

- ❑ Firstly, On the Frequency measures:

```
pca = PCA()  
pca.fit(df_pcl)  
explained_variance_ratio = pca.explained_variance_ratio_  
print(explained_variance_ratio)  
  
[9.64162746e-01 2.33046731e-02 1.11199310e-02 1.41256948e-03  
 7.99702317e-08]
```

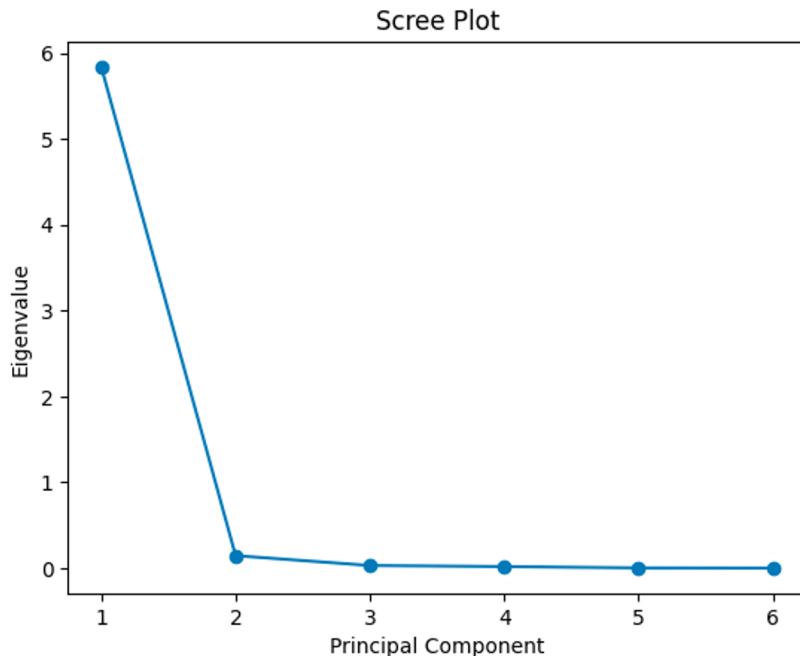
- ❑ As we can see from the scree plot, We can use one principal component to explain the five measures of Jitter.



PCA (2)

- ❑ Repeating the process for the measures of Shimmer.
- ❑ We arrive at the same conclusion
- ❑ One principal component to explain the various measures of Shimmer

```
pca2 = PCA()  
df_pc2 = df_pc[variation_amp]  
df_pc2 = scaler.fit_transform(df_pc2)  
  
pca2.fit(df_pc2)  
explained_variance_ratio = pca2.explained_variance_ratio_  
print(explained_variance_ratio)  
  
[9.67891293e-01 2.39668230e-02 4.97378202e-03 2.88262892e-03  
 2.85467090e-04 6.01341259e-09]
```



PCA Final:

❑ Add these two new principal component columns to the original dataset:

	MDVP:Fo(Hz)	NHR	HNR	RPDE	DFA	spread1	spread2	D2	PPE	status	MDVP:F_Spd	Jitter_pc	Shimmer_pc
0	119.992	0.02211	21.033	0.414783	0.815285	-4.813031	0.266482	2.301442	0.284654	1	82.305	0.934818	1.700278
1	122.400	0.01929	19.085	0.458359	0.819521	-4.075192	0.335590	2.486855	0.368674	1	34.831	1.751692	4.081638
2	116.682	0.01309	20.651	0.429895	0.825288	-4.443179	0.311173	2.342259	0.332634	1	19.556	2.333009	2.865896
3	116.676	0.01353	20.644	0.434969	0.819235	-4.117501	0.334147	2.405554	0.368975	1	26.505	2.020036	3.224675
4	116.014	0.01767	19.649	0.417356	0.823484	-3.747787	0.234513	2.332180	0.410335	1	31.126	3.346368	4.471558
5	120.552	0.01222	21.378	0.415564	0.825069	-4.242867	0.299111	2.187560	0.357775	1	17.375	1.832739	2.152991
6	120.267	0.00607	24.886	0.596040	0.764112	-5.634322	0.257682	1.854785	0.211756	1	22.424	-1.211720	-1.775359

Models (2):

- ❑ The same 4 models are created on the new data.
- ❑ The performance of the models in the new test data:

	Metric	DT	RF	SVM	KNN
0	Accuracy	0.820513	0.846154	0.820513	0.897436
1	F1-Score	0.877193	0.888889	0.877193	0.925926
2	Recall	0.925926	0.888889	0.925926	0.925926
3	Precision	0.833333	0.888889	0.833333	0.925926
4	R2-Score	0.157407	0.277778	0.157407	0.518519

Comparisons:

Model without PCA

	Metric	DT	RF	SVM	KNN
0	Accuracy	0.820513	0.948718	0.846154	0.897436
1	F1-Score	0.881356	0.962963	0.892857	0.928571
2	Recall	0.962963	0.962963	0.925926	0.962963
3	Precision	0.812500	0.962963	0.862069	0.896552
4	R2-Score	0.157407	0.759259	0.277778	0.518519

Model with PCA

	Metric	DT	RF	SVM	KNN
0	Accuracy	0.820513	0.846154	0.820513	0.897436
1	F1-Score	0.877193	0.888889	0.877193	0.925926
2	Recall	0.925926	0.888889	0.925926	0.925926
3	Precision	0.833333	0.888889	0.833333	0.925926
4	R2-Score	0.157407	0.277778	0.157407	0.518519

Feature importance:

- ❑ Feature importance from the Random forest model on the non-pca data.
- ❑ **PPE** is the most important followed by **spread1** and **MDVP:F0(Hz)**

```
# Get feature importances
importances = rfcl.feature_importances_

# Print feature importances
for feature, importance in zip(X.columns, importances):
    print(f"{feature}: {importance}")
```

```
MDVP:F0(Hz): 0.1320668098294375 ←
MDVP:Jitter(Abs): 0.04023009691589492
MDVP:Shimmer: 0.12493439495158931
NHR: 0.06507400866224698
HNR: 0.052912178670557565
RPDE: 0.054599510469434204
DFA: 0.04843645923393191
spread1: 0.13853025374107053 ←
spread2: 0.08964864860571996
D2: 0.05462080274261521
PPE: 0.1443170982033454 ←
MDVP:F_Spd: 0.0546297379741566
```

- ❑ Feature importance from the Random forest model on the pca data.
- ❑ We see the same variables showing up as important.

```
# Get feature importances
importances = rfcl.feature_importances_

# Print feature importances
for feature, importance in zip(X.columns, importances):
    print(f"{feature}: {importance}")
```

```
MDVP:F0(Hz): 0.14001260488596276 ←
NHR: 0.04838287411867583
HNR: 0.045100290961319435
RPDE: 0.045446291652754205
DFA: 0.04897296662768114
spread1: 0.1724148099016464 ←
spread2: 0.08367074620695798
D2: 0.059296514587743716
PPE: 0.1851486311305981 ←
MDVP:F_Spd: 0.05345051076479391
Jitter_pc: 0.045806867271927745
Shimmer_pc: 0.07229689188993882
```