# Comparative Analysis of Machine Learning Models for Parkinson's Disease with PCA

*Jinghua Sun*

## Background

**There** has been an increased interest in speech pattern analysis applications of Parkinsonism for building predictive telediagnosis and telemonitoring models. For this purpose, we have collected a wide variety of voice samples, including sustained vowels, words, and sentences compiled from a set of speaking exercises for people with Parkinson's disease. Main issues in learning from such a dataset that consists of multiple speech recordings per subject is that How predictive these various types, e.g., sustained vowels versus words, of voice samples are in Parkinson's disease (PD) diagnosis?

## Motivation



Parkinson's Disease (PD) is a neurodegenerative disorder affecting millions of people worldwide. Early and accurate diagnosis is crucial for effective management and intervention. Machine learning models have shown promise in aiding the diagnosis of PD based on various clinical features. This study aims to explore the impact of dimensionality reduction, specifically through Principal Component Analysis (PCA), on the performance of different machine learning models for predicting Parkinson's Disease.

Surprisingly, out results indicated a decrease in accuracy after PCA.

## Reference

[1]Sakar, B. E., Isenkul, M. E., Sakar, C. O., Sertbas, A., Gurgen, F., Delil, S., Apaydin, H., & Kursun, O. (2013). Collection and analysis of a Parkinson speech dataset with multiple types of sound recordings. *IEEE Journal of Biomedical and Health Informatics, 17*(4), 828–834.

## Dataset & Machine Learning Models



1. **Jitter features** measure variations in pitch periods

2. **Shimmer features** quantify variations in amplitude and the stability of the voice.

3. **'UPDRS':** Unified Parkinson's Disease Rating Scale
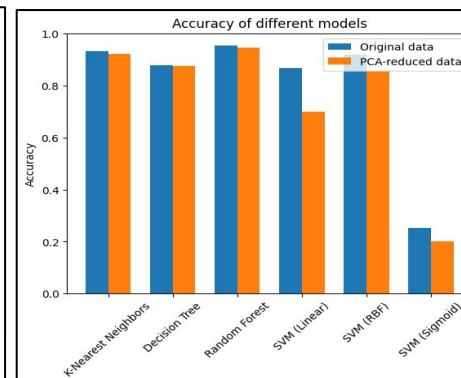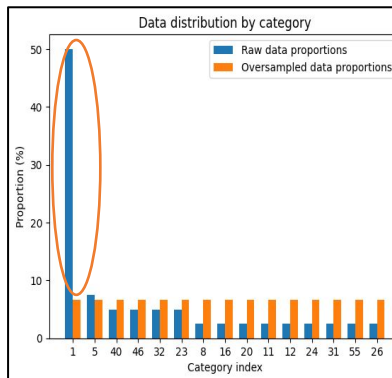
**Principal Component Analysis**

**Random Forest**

```
dfs_pca = []
for gp in gps:
    dfs_pca.append(pca(df[gp]).copy())

for ['jitter_local', 'jitter_local_absolute', 'jitter_rap', 'jitter_ppq5', 'jitter_ddp'], keep 1 features to get an explaine
for ['shimmer_local', 'shimmer_local_db', 'shimmer_apq3', 'shimmer_apq5', 'shimmer_apq11', 'shimmer_data'], keep 1 features
for ['median_pitch', 'mean_pitch', 'standard_dev_pitch', 'min_pitch', 'max_pitch'], keep 2 features to get an explained_vari
for ['AC', 'NTH', 'HTN'], keep 1 features to get an explained_variance_ratio of 0.9464691832031871
for ['num_pulses', 'num_periods', 'mean_period', 'standard_dev_period'], keep 3 features to get an explained_variance_ratio
for ['frac_locally_unvoiced_frames', 'num_voice_breaks', 'degree_of_voice_breaks'], keep 2 features to get an explained_vari
```
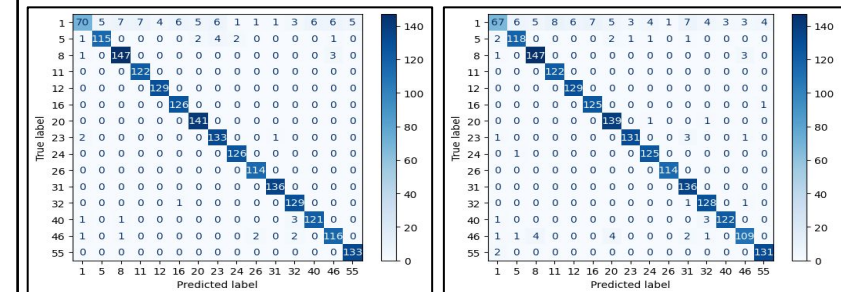
## Analysis Results





Original dataset is imbalanced, the proportion of 1 took nearly **50%** which is much more than other levels
Applied PCA Algo to select the **most important 10 features** which distributed evenly into each category

1. Overall accuracy **decreased** after PCA
2. **Random forest** predict most accurately both with and without PCA
3. **SVM(RBF)** predicted second most accurately both with and without PCA
4. The least accuracy is predicted by **SVM(SIG)** both with and without PCA
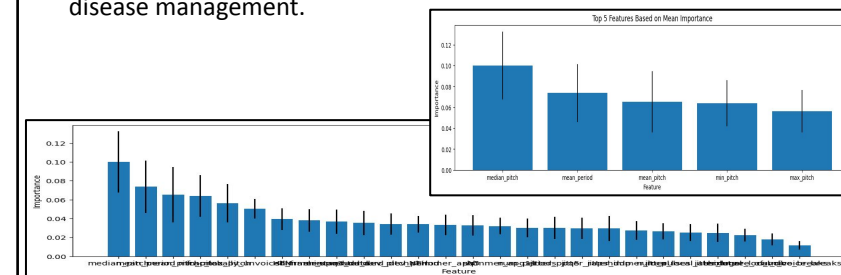
## Interpretation of Results

### Random Forest



### TOP 5 Importance Features

**MEDIAN_PITCH,MEAN_PERIOD,MEAN_PITCH,MIN_PITCH,MAX_PITCH**

- **MEDIAN_PITCH :** Could be linked to vocal fold dynamics, muscle control, or other factors associated with Parkinson's disease.
- **MEAN_PERIOD :** Related to voice quality and potential motor control issues associated with Parkinson's disease.
- **Healthcare Professionals:** Decision-making & Treatment planning.
- **Researchers:** Speech as a biomarker for neurological conditions.
- **Patients:** A convenient and regular monitoring method, improving disease management.



### SVM-RBF