

Introduction

Photoplethysmography (PPG) is a non-invasive optical technique that has been known to be a successful indicator of Blood Pressure (BP). In this study, we developed classical and deep learning models on PPG signal data in order to predict Diastolic and Systolic BP. We constructed several regression and classification models and implemented feature reduction algorithms to model BP with fewer features than standard models. Our models (trained on exclusively PPG data) performed at similar or better accuracy levels than recent, state-of-the-art models. Further, we showed that it is possible to obtain comparable results with significantly less features.

Dataset

We used the "Cuff-Less Blood Pressure Estimation Data Set" from UCI Machine Learning Repository which contains processed data from the MIMIC II database. There are 12000 instances of PPG, ABP, and ECG signals; the ABP dataset provided the ground truth for Systolic Blood Pressure (SBP) and Diastolic Blood Pressure (DBP). After filtering, 9488 examples were kept. The SBP and DBP values were divided into 5 classes as recommended by American Heart Association for the classification algorithms.

Classes	SBP	and/or	DBP
Low Blood Pressure	<90	or	<60
Normal Blood Pressure	<120	and	<80
Elevated Blood Pressure	120 - 129	and	<80
Hypertension Stage I	130-139	or	80-89
Hypertension Stage II	>140	or	>89

Figure 1. Blood Pressure Classes

Feature Extraction

Since our goal is to predict SBP and DBP solely using the PPG features, we only extracted features from PPG waveforms. We extracted features from PPG waveforms using techniques devised by Thambiraj et al [2] and Kachuee et al [1]. A total of 8 physiological features, and 20 temporal features were extracted.

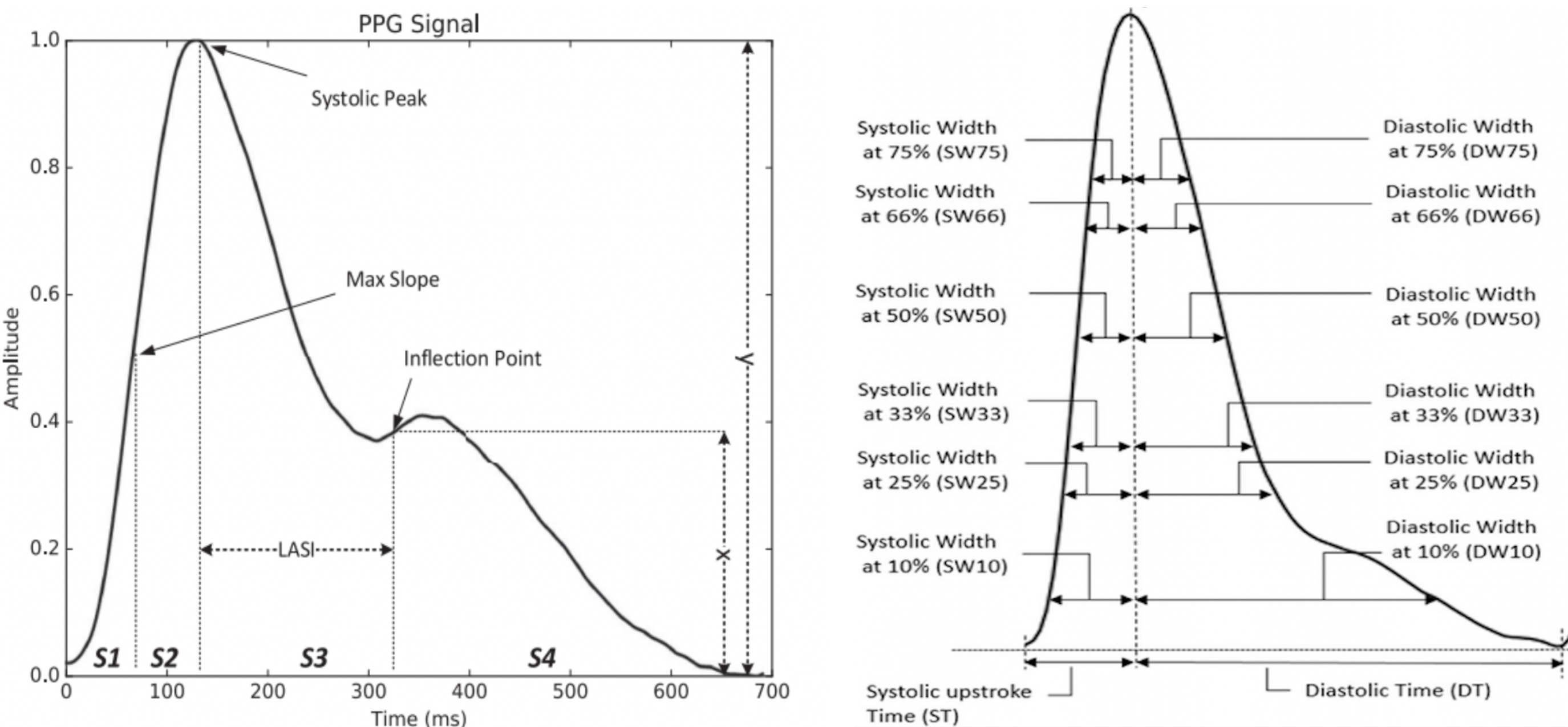


Figure 2. Physiological & Temporal Features

Models

- Baseline Models** - For BP *value* prediction, we used Linear Regression, Ridge Regression, Support Vector Regression, and LSTM as our baseline. For BP Classification, we used Softmax as baseline.
- Neural Network** - Our Neural Network architecture is as shown in Figure 4. A learning rate of 0.001 was employed, with training over 1000 epochs. MSE and Cross Entropy were used as loss for regression and classification, respectively.

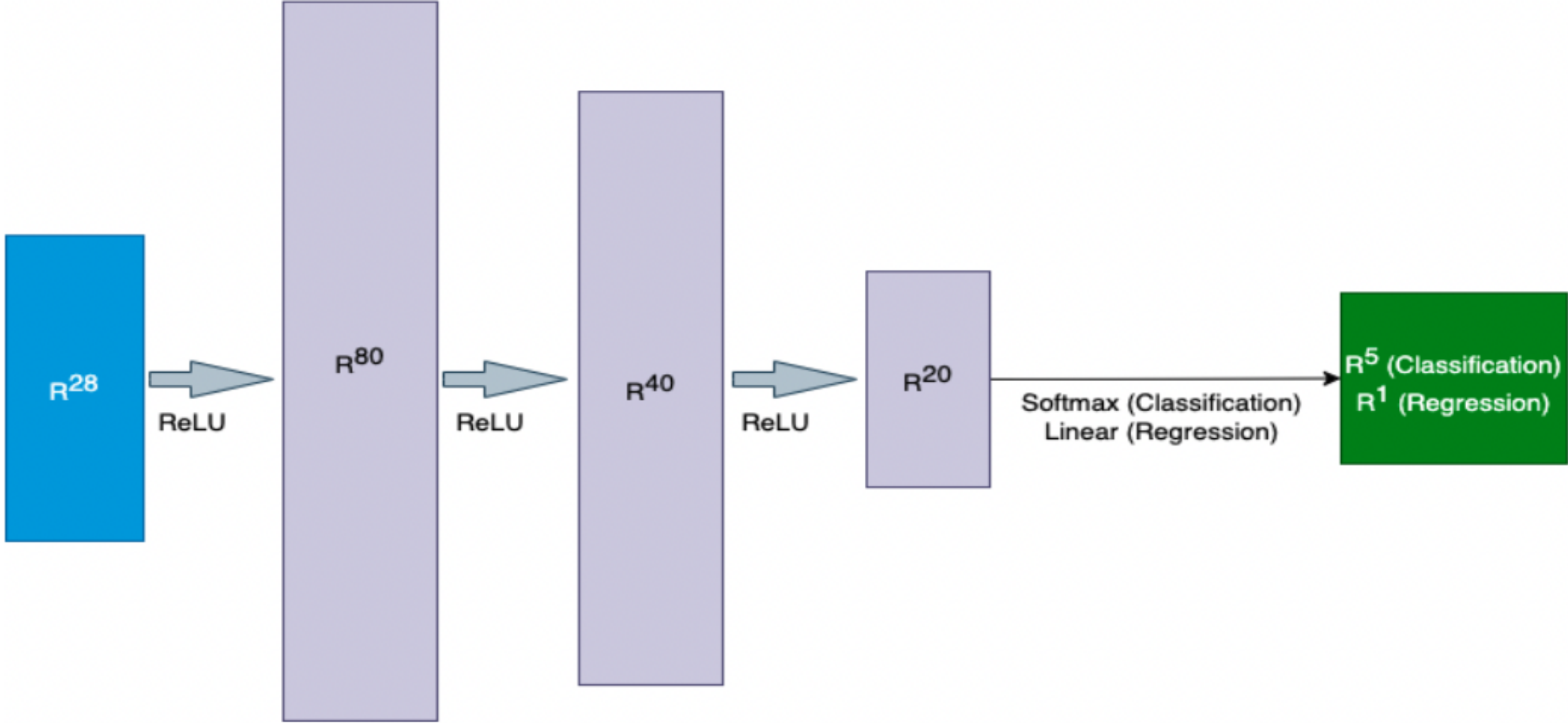


Figure 3. Neural Network Architecture (Regression & Classification)

- Random Forest** - Random forest (RF) regression is a type of ensemble regression algorithm that works by training multiple decision trees on a dataset and then combining the predictions of each tree to make a final prediction. Random forest is an effective method for dealing with nonlinear and complex relationships in the data. Our random forest models use 100 trees for both regression and classification.

Feature Selection

We employed multiple feature selection techniques for removing unnecessary or redundant features as well as for identifying the features that are most indicative of blood pressure. The following three methods performed better than all other techniques.

- Pearson’s method** - Pearson’s correlation coefficient is a statistical measure that is used to determine the strength and direction of the linear relationship between two variables, defined by:

$$\rho_{X,Y} = \frac{\mathbb{E}[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

- One-way ANOVA** - One-way ANOVA (short for "analysis of variance") tests for statistically significant differences in the means between two or more independent groups, characterized by:

$$F = \frac{\sum_{i=1}^K n_i (\bar{Y}_i - \bar{Y})^2 / (K - 1)}{\sum_{ij=1}^n (Y_{ij} - \bar{Y}_i)^2 / (N - K)}$$

- Random Forest Model Selection** For our Random Forest model, we employed a feature selection technique provided from scikit-learn known as SelectFromModel. For random forest, the importance of a feature is estimated by the mean decrease in impurity (MDI).

$$MDI = \frac{1}{n \text{ tree}} \sum_{t=1}^{n \text{ tree}} (MP_{tj} - M_{tj})$$

Results

We used 6641 training and 2847 testing examples. Performance of feature-reduced models are found in our paper. For example, with 16 features, the test accuracy is 99% for random forest and 96% for NNs.

Model	Task	MAE for DBP		MAE for SBP		Classification Accuracy (%)	
		Train	Test	Train	Test	Train	Test
Linear Regression	Regression	15.97	16.08	16.09	16.01	-	-
Ridge Regression	Regression	17.07	16.98	16.63	16.18	-	-
SVR	Regression	17.42	17.19	17.42	17.19	-	-
LSTM	Regression	7.393	7.543	7.100	7.293	-	-
Random Forest	Regression	1.018	1.602	1.344	2.082	-	-
Neural Network	Regression	2.565	3.139	1.826	4.767	-	-
Softmax Regression	Classification	-	-	-	-	62	58
Random Forest	Classification	-	-	-	-	100	99
Neural Network	Classification	-	-	-	-	99	96

Figure 4. Performance of all models for both regression and classification, trained with all 28 features

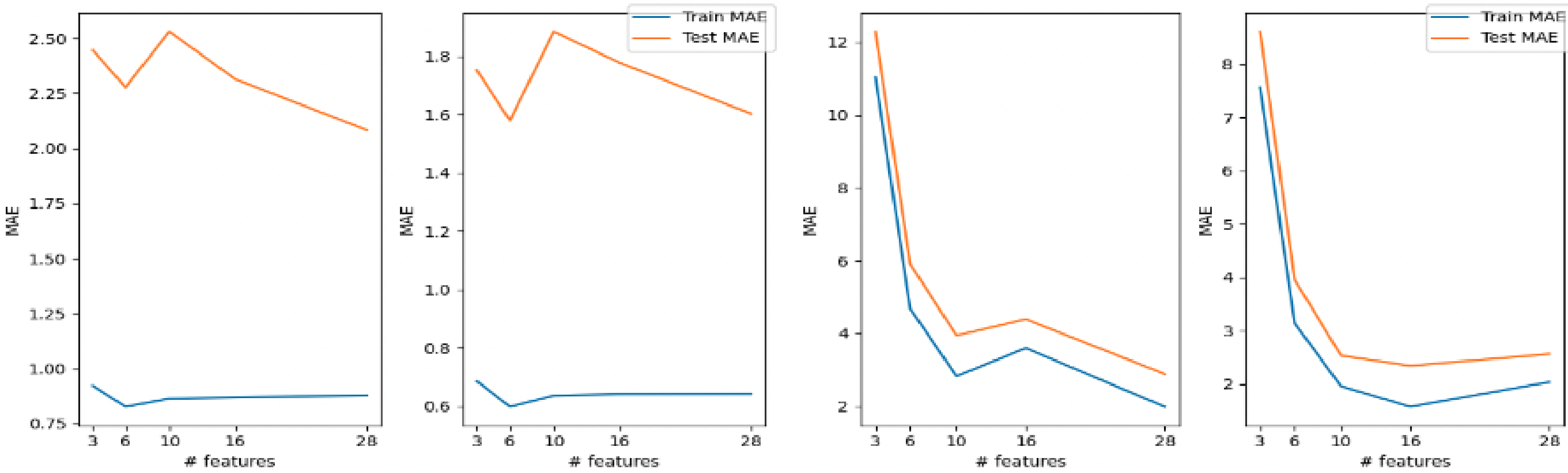


Figure 5. MAE of SBP, DBP after Feature Reduction for Random Forest & Neural Network

Discussion

Our results show that our selected features are able to predict systolic and diastolic blood pressure more accurately than most existing models in the literature. Our Random Forest Regressor performs best for BP prediction, while both Random Forest and Feed-Forward Neural Networks perform strongly for BP Classification. We also showed that we can recreate similar high accuracy by selecting significantly fewer features. Our results provide the first analysis of specific PPG features that contribute to BP prediction.

Future Work

- Implement a Deep Multitask Network to optimize loss with respect to correlated features
- Implement CNN-BiLSTM network to extract temporal features and predict the ABP pulse wave itself

References

[1] Mohammad Kachuee, Mohammad Mahdi Kiani, Hoda Mohammadzade, and Mahdi Shabany. Cuffless blood pressure estimation algorithms for continuous health-care monitoring. *IEEE Transactions on Biomedical Engineering*, 64(4):859–869, 2017.

[2] Geerthy Thambiraj, Uma Gandhi, Umopathy Mangalanathan, V Jeya Maria Jose, and M Anand. Investigation on the effect of womersley number, ecg and ppg features for cuff less blood pressure estimation using machine learning. *Biomedical Signal Processing and Control*, 60:101942, 2020.