

Fusing Features based on Signal Properties and TimeNet for Time Series Classification

Arijit Ukil, Pankaj Malhotra, Soma Bandyopadhyay, Tulika Bose, Ishan Sahu,
Ayan Mukherjee, Lovekesh Vig, Arpan Pal, and Gautam Shroff
{arijit.ukil,malhotra.pankaj,soma.bandyopadhyay,tulika.bose,ishan.sahu,
ayan.m4, lovekesh.vig,arpan.pal,gautam.shroff}@tcs.com

TCS Research, India

Abstract. Automated feature extraction from time series to capture statistical, temporal, spectral, and morphological properties is highly desirable but challenging due to diverse nature of real-world time series applications. In this paper, we consider extracting a rich and robust set of time series features encompassing signal processing based features as well as generic hierarchical features extracted via deep neural networks. We present *SPGF-TimeNet*: a generic feature extractor for time series that allows fusion of signal processing, information-theoretic, and statistical features (Signal Properties based Generic Features (*SPGF*)) with features from an off-the-shelf pre-trained deep recurrent neural network (*TimeNet*). Through empirical evaluation on diverse benchmark datasets from the UCR Time Series Classification (TSC) Archive, we show that classifiers trained on *SPGF-TimeNet*-based hybrid and generic features outperform state-of-the-art TSC algorithms such as BOSS, while being computationally efficient.

1 Introduction

Extracting custom features for any Internet of Things (IoT) application often requires costly, and sometimes impractical, expert intervention and domain knowledge. There is an increasing need to involve little-to-no expert intervention in these applications. Representation of time series in terms of a rich set of generic features can address this requirement and potentially automate time series analysis. Deep neural networks pre-trained on several diverse time series have been found to provide useful generic hierarchical features for unseen datasets, e.g. TimeNet based on multilayered recurrent neural network (RNN) [1, 2] and Universal Encoder based on convolutional neural networks [3]. On the other hand, Signal Properties based Generic Features (SPGF) have also been found to yield promising results in automated feature extraction for time series [4]¹.

Here, we consider an effective combination of features using deep learning (TimeNet) and traditional signal processing (SPGF), and propose *SPGF-TimeNet* as a generic time series feature extractor. *SPGF-TimeNet* is completely automatable due to unsupervised feature extraction, and therefore, invariably minimizes human effort, bias and computational cost (complexity of feature extraction using *SPGF-TimeNet* varies linearly with length L of time

¹Winning entry of CiNC challenge 2017 (<https://physionet.org/challenge/2017/>)

series $\approx \mathcal{O}(L)$) while potentially maximizing scalability due to the absence of human-labored feature engineering. We consider UCR TSC datasets [5] for empirical evaluation and observe significantly better classification performance of *SPGF-TimeNet* over state-of-the-art TSC algorithms like BOSS [6].

2 SPGF-TimeNet based Time Series Classification

Consider a univariate time series signal $x_{1...L} = x_1, x_2, \dots, x_L$ ($x_i \in \mathbb{R}$) of length L . Let $\mathbf{z}_S \in \mathbb{R}^{n_1}$ represent the n_1 features extracted via SPGF, and $\mathbf{z}_T \in \mathbb{R}^{n_2}$ represent the n_2 features extracted via TimeNet (TN). The final $(n_1 + n_2)$ -dimensional feature vector representation \mathbf{z}_{ST} of $x_{1...L}$ is then given by concatenation of features from SPGF and TN as $\mathbf{z}_{ST} = [\mathbf{z}_S, \mathbf{z}_T] \in \mathbb{R}^{n_1+n_2}$. We next explain the construction of \mathbf{z}_S and \mathbf{z}_T from $x_{1...L}$, and then describe how these are combined for the task of TSC via a classifier such as a non-linear SVM [7].

2.1 SPGF-based Feature Extraction (\mathbf{z}_S)

We consider extracting a diverse set of features to capture the intrinsic morphological and statistical properties as well as inherent randomness and regularity in the time series. To this end, we extract features from three domains [8, 9, 10]: 1) *temporal* (\mathcal{T}): original time series, 2) *spectral* (\mathcal{S}): transforming the time series into frequency domain (Fast Fourier Transform), and 3) *wavelet* (\mathcal{W}): capturing frequency variations over time. We apply Discrete Wavelet Transform (DWT) upto 4th level with Daubechies wavelets 4 (db4) as mother wavelet in \mathcal{W} domain as the corresponding coefficients capture maximum amount of signal energy [11]. In order to capture the global and local properties of a time series, we extract macro and micro-level features from representations of time series \mathcal{T} , \mathcal{S} and \mathcal{W} domains. More specifically, the macro-level features are extracted by considering the entire time series $x_{1...L}$ in \mathcal{T} domain or the corresponding transformations in \mathcal{S} and \mathcal{W} domains at once, while the micro-level features are extracted from every *window* $x_{t+1...t+\tau}$ of time series in \mathcal{T} domain or corresponding transformations of the windows in \mathcal{S} and \mathcal{W} domains. We obtain n_ω non-overlapping *windows* of time series $x_{1...L}$ with each window $x_{t+1...t+\tau}$ of length $\tau \approx \frac{L}{n_\omega}$, with $t = 0, \tau, 2\tau, \dots, (n_\omega - 1)\tau$.

We extract several statistical and information-theoretic features as depicted in Fig. 1(a) corresponding to the macro and micro-level features. Examples of statistical features considered include mean, standard deviation, kurtosis, skewness, Box Pierce statistics and Hurst exponent [12]. Examples of information-theoretic features include Shannon, Tsallis, and Renyi entropies [13]. *Shannon_entropy*($x_{1...L}$) is a macro-level feature which represents the entropy of the entire time series. Any micro-level feature calculated for each of the n_ω windows eventually results in two features given by the mean and standard deviation of the n_ω values of the feature (refer Fig. 1(a)). For example, one of the micro-level features *mean of windowed kurtosis* in temporal domain is given by the mean of kurtosis values of the n_ω time series windows, representing the degree of tailedness. An exemplary micro-level feature is described in Section 3.

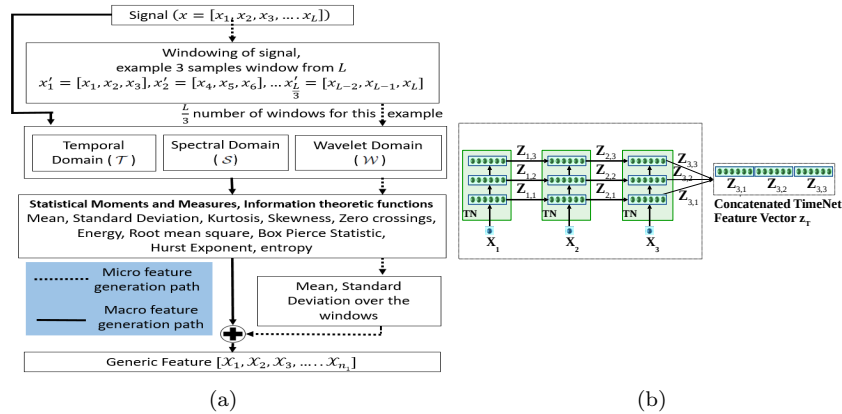


Fig. 1: (a) SPGF Feature Extraction Details, (b) TimeNet-based Feature Extraction. TimeNet (TN) is shown unrolled for $L = 3$. (Best visible when zoomed.)

2.2 TimeNet-based Feature Extraction (z_T)

We consider extracting hierarchical features from TimeNet [1, 2]. TimeNet is a pre-trained off-the-shelf feature extractor for univariate time series. It consists of three hidden layers with 60 Gated Recurrent Units (GRUs) each. The univariate input time series is mapped by TimeNet to 180-dimensional feature vector such that $n_2 = 180$, where each dimension corresponds to final output of one of the 60 GRUs in the 3 recurrent layers. TimeNet has been shown to be effective for diverse time series of varying length ($L \leq 512$) across diverse domains.

TimeNet is the encoder part of the autoencoder consisting of an encoder RNN f_E and a decoder RNN f_D trained simultaneously on 24 diverse time series datasets using unsupervised sequence-to-sequence learning framework as shown in Figure 1(b). The parameters \mathbf{W}_E of the encoder RNN f_E is obtained by training the autoencoder via reconstruction task so that for input $x_{1...L}$, the target output time series $x_{L...1} = x_L, x_{L-1}, \dots, x_1$ is reverse of the input. The RNN encoder f_E maps the univariate input time-series $x_{1...L}$ non-linearly to a fixed-dimensional feature representation \mathbf{z}_L at the L -th time step: $\mathbf{z}_L = f_E(x_{1...L}; \mathbf{W}_E)$. The feature vector \mathbf{z}_L is a concatenation of the hidden states $\mathbf{z}_{L,l}$ ($l = 1, 2, 3$) from the three layers. During training, this is followed by a non-linear mapping of \mathbf{z}_L to univariate time series: $\hat{x}_{L...1} = f_D(\mathbf{z}_L; \mathbf{W}_D)$ via an RNN decoder f_D ; where \mathbf{W}_E and \mathbf{W}_D are the parameters of the encoder and decoder, respectively. The mean squared reconstruction error is used as loss function to jointly train the encoder (TimeNet) and decoder. Since the decoder relies on \mathbf{z}_L as the only input to reconstruct the time series, the encoder gets trained to capture all the relevant information in the time series $x_{1...L}$ into the fixed-dimensional vector \mathbf{z}_L . We, therefore, use $\mathbf{z}_L \equiv \mathbf{z}_T \in \mathbb{R}^{180}$, i.e. use the final hidden state of TimeNet \mathbf{z}_L after processing $x_{1...L}$ as the feature vector extracted via TimeNet.

2.3 Using feature vector for TSC

Consider a labeled training set with N instances $\{x_{1\dots L}^i, y^i\}_{i=1}^N$ of univariate time series $x_{1\dots L}^i$ and corresponding class label y^i . The $(n_1 + n_2)$ -dimensional feature vector representation \mathbf{z}_{ST}^i of $x_{1\dots L}^i$ is obtained as $\mathbf{z}_{ST}^i = [\mathbf{z}_S^i, \mathbf{z}_T^i]$. Thus, we convert the training set to *SPGF-TN* feature space to obtain $\{\mathbf{z}_{ST}^i, y^i\}_{i=1}^N$. The \mathbf{z}_{ST} features are appropriately z-normalized and used to train an SVM-based classifier with Radial Basis Function (RBF) kernel (similar to [1]). The hyperparameters γ (kernel coefficient) and ν (rejection rate) of SVM-RBF are obtained using 5-fold cross-validated grid search on the training set over the logarithmic grid of both γ and ν in range 10^{-3} to 10^3 .

3 Experiments and Analysis

We consider a diverse subset of the UCR TSC Archive to test the generality of the *SPGF-TN* features and corresponding classifier *SPGF-TN-C* using same setup and data splits as in [1], and consider DTW-C [5], BOSS [6] and TimeNet (TN-C) [1] as baselines for comparison. We consider number of windows $n_\omega = 10$, $n_1 = 392$ and $n_2 = 180$ such that we obtain 572-dimensional feature vector via *SPGF-TN*, i.e. $\mathbf{z}_{ST} \in \mathbb{R}^{572}$. We obtain a total $n_1 = 392$ *SPGF* features: 54 in temporal, 48 in frequency, and 290 in wavelet domains. Overall, *SPGF* comprises 248 micro-level and 144 macro-level features. For instance, if $n_\omega = 3$ and kurtosis values for the three windows are 19.09, 4.75 and 3.31 respectively, the value of the micro-level feature *mean of windowed kurtosis* in temporal domain is 9.05 and that of *standard deviation of windowed kurtosis* is 8.72. The error rates given by $1 - \text{accuracy}$ are summarized in Table 1. The win/tie counts depict the number of data sets in which the respective model individually performs better or is equivalent to the best of DTW-C and BOSS. We observe that *SPGF-TN-C* outperforms both DTW-C and BOSS in 22/30 of the cases. We further observe that *SPGF-TN-C* is better than either of *SPGF-C* or *TN-C*. This proves the advantage of combining deep learning based features (TN) and signal processing based features (SPGF), and the richness of the feature space of *SPGF-TN*. In Fig. 2 (a),(b), we further illustrate this by mapping the 500-dimensional ($L = 500$) original time series test instances and the corresponding 572-dimensional feature vectors of FordB dataset to 2-D space using t-SNE (blue and green colors depict class labels): we find a noticeable separation in the two classes in the *SPGF-TN* space (Fig. 2(b)) while no separation in the original space (Fig. 2(a)).

Further, we note that sequential processing in RNNs and various transformations in *SPGF* are linear in the length L of time series. Therefore, the inference cost of *SPGF-TN* is linear w.r.t. time series length $\approx \mathcal{O}(L)$. Fig. 2(c) depicts this behavior in terms of execution times of *SPGF*, *TimeNet*, and *SPGF-TN* over test time series of varying lengths². This is of high practical importance as algorithmic complexity can play a major role in deciding the selection of appropriate TSC technique [5]. For instance, the state-of-the-art TSC algorithm

²Experiments done on system with Intel Xeon processor with 2.60 GHz and 128 GB RAM

COTE takes ensemble of 35 classifiers [5], and is computationally very expensive. Similarly, BOSS [6] has quadratic inference complexity in length of time series.

Table 1: Comparison of Classification Error Rates on UCR TSC Datasets.

Dataset	L	DTW-C [5]	BOSS [6]	SPGF-TN-C (ours)	TN-C [1]	SPGF-C
Synthetic Control	60	0.017	0.033	0.017	0.013	0.023
PhalangesOC	80	0.239	0.228	0.184	0.207	0.187
DistalPhalanxOAG	80	0.228	0.252	0.163	0.223	0.173
DistalPhalanxOC	80	0.232	0.272	0.155	0.188	0.157
DistalPhalanxTW	80	0.272	0.324	0.218	0.208	0.223
MiddlePhalanxOAG	80	0.253	0.455	0.220	0.210	0.228
MiddlePhalanxOC	80	0.318	0.220	0.330	0.270	0.355
MiddlePhalanxTW	80	0.419	0.455	0.363	0.363	0.371
ProximalPhalanxOAG	80	0.215	0.166	0.151	0.146	0.151
ProximalPhalanxOC	80	0.210	0.151	0.124	0.175	0.134
ProximalPhalanxTW	80	0.263	0.200	0.193	0.195	0.200
ElectricDevices	96	0.376	0.201	0.275	0.267	0.306
MedicalImages	99	0.253	0.282	0.241	0.250	0.258
Swedish Leaf	128	0.157	0.078	0.064	0.102	0.064
Two Patterns	128	0.002	0.007	0.000	0.000	0.221
ECG5000	140	0.075	0.059	0.058	0.069	0.060
ECGFiveDays	136	0.203	0.000	0.051	0.074	0.047
Wafer	152	0.005	0.005	0.000	0.005	0.000
ChlorineConcentration	166	0.350	0.339	0.279	0.269	0.306
Adiac	176	0.391	0.235	0.212	0.322	0.246
Strawberry	235	0.062	0.024	0.049	0.062	0.047
Cricket_X	300	0.236	0.264	0.264	0.300	0.351
Cricket_Y	300	0.197	0.246	0.280	0.338	0.387
Cricket_Z	300	0.180	0.254	0.231	0.308	0.318
uWaveGestureLib_X	315	0.227	0.238	0.171	0.214	0.256
uWaveGestureLib_Y	315	0.301	0.315	0.237	0.311	0.321
uWaveGestureLib_Z	315	0.322	0.305	0.226	0.281	0.286
Yoga	426	0.155	0.082	0.145	0.160	0.177
FordA	500	0.341	0.07	0.068	0.219	0.058
FordB	500	0.414	0.289	0.108	0.263	0.158
Wins or ties over both DTW-C & BOSS	-	-	-	22/30	16/30	15/30

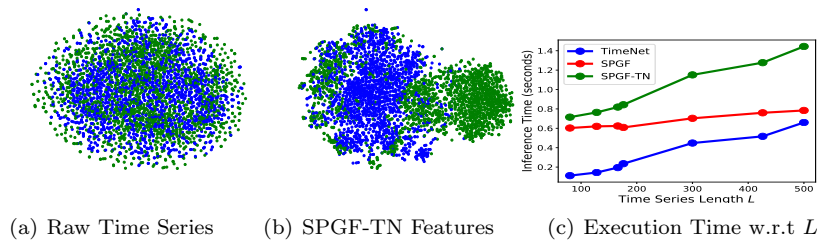


Fig. 2: t-SNE scatter plot for FordB test dataset: (a) Raw Time series (b) *SPGF-TN* features; (c) Execution times of *SPGF*, *TimeNet*, and *SPGF-TimeNet*.

4 Conclusion

We have proposed *SPGF-TimeNet* for extracting generic features from time series. *SPGF-TimeNet* combines the advantage of signal processing and deep

learning to yield generic feature set that are observed to be useful for diverse time series classification (TSC) tasks with variable-length time series in a domain-agnostic way. It ensures effective and efficient learning particularly due to its richness of the feature space. It outperforms existing state-of-the-art in more than 73% of datasets considered. In future, it will be interesting to evaluate SPGF-TimeNet for multivariate TSC tasks (e.g. as in [2]). We also plan to further augment the feature space of SPGF using graph signal processing and dictionary learning techniques, and also train a bigger TimeNet to enhance the feature space of *SPGF-TimeNet*. The sequential processing in RNNs makes it expensive to train TimeNet for long time series - we plan to exploit convolution neural networks based TimeNet that is faster to train (e.g. as in [3]).

References

- [1] Pankaj Malhotra, Vishnu TV, Lovekesh Vig, Puneet Agarwal, and Gautam Shroff. Timenet: Pre-trained deep recurrent neural network for time series classification. In *25th European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning*, pages 607–612, 2017.
- [2] Priyanka Gupta, Pankaj Malhotra, Lovekesh Vig, and Gautam Shroff. Using features from pre-trained timenet for clinical predictions. In *The 3rd International Workshop on Knowledge Discovery in Healthcare Data at IJCAI*, 2018.
- [3] Joan Serra, Santiago Pascual, and Alexandros Karatzoglou. Towards a universal neural network encoder for time series. *arXiv preprint arXiv:1805.03908*, 2018.
- [4] S. Datta, C. Puri, A. Mukherjee, R. Banerjee, A. D. Choudhury, R. Singh, A. Ukil, S. Bandyopadhyay, A. Pal, and S. Khandelwal. Identifying normal, af and other abnormal ecg rhythms using a cascaded binary classifier. In *2017 Computing in Cardiology (CinC)*, pages 1–4, Sept 2017.
- [5] Anthony Bagnall, Jason Lines, Aaron Bostrom, James Large, and Eamonn Keogh. The great time series classification bake off: A review and experimental evaluation of recent algorithmic advances. *Data Min. Knowl. Discov.*, 31(3):606–660, May 2017.
- [6] Patrick Schäfer. The BOSS is concerned with time series classification in the presence of noise. *Data Min. Knowl. Discov.*, 29(6):1505–1530, November 2015.
- [7] Christopher JC Burges. A tutorial on support vector machines for pattern recognition. *Data mining and knowledge discovery*, 2(2):121–167, 1998.
- [8] George EP Box, Gwilym M Jenkins, Gregory C Reinsel, and Greta M Ljung. *Time series analysis: forecasting and control*. John Wiley & Sons, 2015.
- [9] Martin Vetterli, Jelena Kovačević, and Vivek K Goyal. *Foundations of signal processing*. Cambridge University Press, 2014.
- [10] S. Banerjee, T. Chattopadhyay, A. Pal, and U. Garain. Automation of feature engineering for iot analytics. *ACM SIGBED Review*, pages 24 – 30, 2018.
- [11] Nishchal K Verma, Rahul Kumar Sevakula, Sonal Dixit, and Al Salour. Intelligent condition based monitoring using acoustic signals for air compressors. *IEEE Transactions on Reliability*, 65(1):291–309, 2016.
- [12] Harold Edwin Hurst. Long-term storage capacity of reservoirs. *Trans. Amer. Soc. Civil Eng.*, 116:770–799, 1951.
- [13] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.