

Hybrid Model for Interpretable Time Series Analysis

A THESIS

Presented to the Department of Computer Science and Computer Engineering

California State University, Long Beach

In Partial Fulfillment

of the Requirements for the Degree

Master of Science in Computer Science

Option in Computer Science

Committee Members:

Ju Cheol Moon, Ph.D

Roman Tankelevich, Ph.D

Seok-Chul Kwon, Ph.D

College Designee:

Hamid Rahai, Ph.D.

By Ruben Rosales

B.S., 2016, California State University, Fullerton

May 2020

Abstract

In this thesis, we study interpretable machine learning as applied to complex-valued time-series data. Scientists have studied the use of several machine learning methods such as Convolutional Neural Networks, Recurrent Neural Networks, and Support Vector Machines for time series classification. These methods, however, fall short of allowing users to visualize patterns, which can be a necessity for adoption in highly scrutinized industries.

To address this issue, we propose an interpretable hybrid model that can be extended to any time series dataset. In the majority of existing work in the field of interpretability, black-box models tend to outperform white-box models consistently. However, instead of relying purely on one method, we propose a collaboration between the two. This allows for the performance of a black-box model while providing a more precise visualization of our dataset.

To accomplish this, we created a framework composed of two models, an ensemble method of classifiers that functions as a black-box model, and a white-box model that encodes time series as images. The white-box model attempts to classify them through a neural network and outputs all images correctly classified in both black-box and white-box models.

Acknowledgments

I want to thank Dr. Moon for his guidance and for allowing me to work on a project this extensive. I would also like to thank my friends and family for their continuous support throughout this thesis.

Contents

Abstract		ii
Acknowledgments		iii
Figures		v
Chapter 1	Introduction	1
Chapter 2	Related Work	2
Chapter 3	Data	3
Chapter 4	Preprocessing	6
Chapter 5	Models	12
Chapter 6	Markov Transition Fields	18
Chapter 7	Main Results	21
Chapter 8	Conclusion	24

Figures

1.1	Overview of our proposed hybrid model.	1
3.1	Size of each dataset.	3
3.2	Class representation of dataset A and C.	4
3.3	Class representation of dataset B.	5
4.1	Example of short-time fourier transform formula [1].	7
4.2	Example of Short-Time Fourier Transform.	8
4.3	Example of a wavelet decomposition network.	11
5.1	Example of a Convolutional Neural Network [2].	13
5.2	Example of a LSTM unit [5].	14
5.3	CNN architecture used in our black-box model.	15
5.4	LSTM architecture used in our black-box model.	16
5.5	Overview of our black-box model.	17
6.1	Example of QxQ quantile bin matrix used to calculate MTF [18].	18
6.2	Markov Transition Field where W_{ij} denotes the transition probability from quantile i to j	19
6.3	CNN architecture for classifying images generated through MTF.	20
7.1	Accuracy on test data for our black-box model.	22
7.2	Examples of images generated through MTF.	23

Chapter 1

Introduction

As time-series datasets become ever more popular, the need for deep learning models becomes that much greater. In our black-box model, we analyze deep learning models such as Convolutional Neural Networks, which have become increasingly popular and widely used in the last decade [4]. However, their complex nature results in a lack of understanding of its inner workings [17]. As the use and complexity of black-box models rises, so too does the need for an interpretable accessory. For this reason, we propose a hybrid model consisting of a black-box for classification and an interpretable white-box which highlights characteristics of a time series for users to understand.

The goal of our work is not to find an alternative to black-box models, but to provide a way for users to have an understanding of what features are prevalent in their dataset.

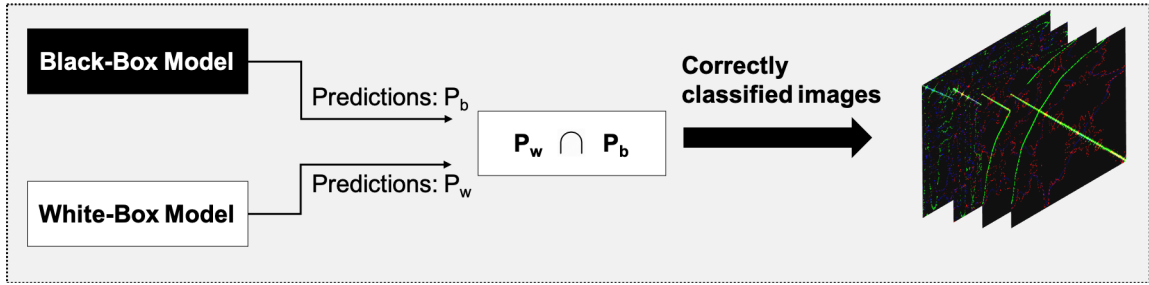


Figure 1.1: Overview of our proposed hybrid model.

Chapter 2

Related Work

The study of interpretable hybrid models has been worked on extensively in the past few years as deep learning has become increasingly popular. Scientists take different approaches as to what makes a hybrid model interpretable and apply it to different domains, such as computer vision or time series. For example, Wang [16], focuses on creating a black-box model that utilizes wavelet decomposition for the classification of time-series that outperforms similar architectures. However, their white-box model focuses on analyzing the importance of their model's layers. This helps users identify the importance of wavelet decomposition but not necessarily understand their dataset.

Additionally, other work we found tends to focus on non-time-series data or creating a general framework defining the trade-off between interpretability and black-box methods. For instance, [17] defines a set of rules to achieve a high level of classification with interpretability but does not give any competitive methods to achieve this.

In order to fill the gap for a high performing interpretable model for time-series classification, we propose our hybrid model, which focuses on extensibility and robustness. We achieve this through a combination of deep learning models for classification as well as a unique white-box method to explain features within classes of a given dataset.

Chapter 3

Data

We analyzed our model on three radio signal datasets, which we will refer to as dataset A, B, and C throughout this thesis. These datasets were generated through Spatial Modulation, which is a transmission technique that uses multiple antennae to map a block of information bits to two units, a symbol chosen from a constellation diagram, and a unique transmit antenna number that is chosen from a set of transmit antennae [11]. The method in which transmit antennae send and receive radio waves can be described through polarization. Two popular basis polarizations are horizontal linear, H, and vertical linear, V [13]. Our datasets contain data horizontally transmitted and vertically received (HV), as well as data vertically transmitted and vertically received (VV).

Dataset	A	B	C
Number of signals	1000	1000	10,000

Figure 3.1: Size of each dataset.

The two properties, HV and VV, are given to us in a raw and discrete format. Both properties consist of time and amplitude values but vary in length and in how they are processed. The raw signal consists of over 31,000 data points with the exact time it was received. The size of the discrete data can vary from 22 to 44 data points and is a processed representation of the raw signal that has equally spaced values. In this thesis, we focus on discrete data given that it is closer to datasets in real-world situations.

Dataset Classes

Dataset A and C were generated through similar methods and can be visually seen in Figure 3.2. There are three base stations, in which each base station corresponds to a class, sending information to a single receiver antenna.

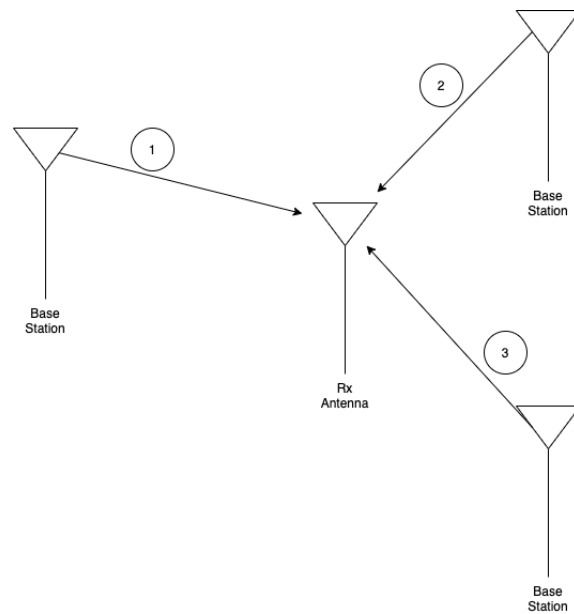


Figure 3.2: Class representation of dataset A and C.

Dataset B differs from dataset A and C in how it transmits data. In 3.3, we can see that there are three different antennas at the same station sending information to a single receiver station. The problem of classification becomes more challenging because they are very close together.

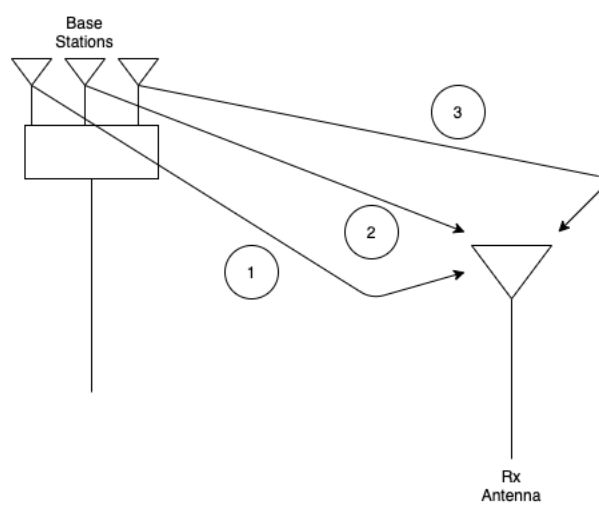


Figure 3.3: Class representation of dataset B.

Chapter 4

Preprocessing

The radio signal data is presented in a complex-valued format that is unusable in a typical neural network, so in an attempt to make our model robust and reduce any complexity, we focused on using real values as features.

We tried numerous methods to extract real-valued features, including Fourier Transform, Short-Time Fourier Transform, a custom sliding window method, and polar coordinates taken directly from the signal's amplitude/phase value. Before applying any of these methods, we normalized all data using L2 normalization.

Transformations

Fourier Transform

The Fourier Transform (FT) is a mathematical tool that decomposes any function into a sum of sinusoidal basis functions. Each basis function is a complex value of a frequency, allowing us to view our data in the frequency domain as opposed to the amplitude domain. Our dataset are discrete in time domain so we focus on discrete fourier transforms. The discrete fourier transform of a signal x can be defined as

$$X_k = \sum_{n=0}^{N-1} x(n) e^{-\frac{i2\pi kn}{N}} \quad (4.1)$$

where $x(t_n)$ is the input signal at time n , $X(k)$ is the k th spectral sample [14].

One of the shortcomings of the FT is rooted in the Heisenberg Uncertainty Principle (HUP) [7]. The HUP states that the position and velocity of an object cannot be simultaneously measured, which can be applied to the time-frequency information of a signal. This means we cannot know which spectral components exist at any given time. The closest we can get is sampling at different ranges of time and finding a range of frequencies within that time frame. This method is described as the Short-Time Fourier Transform.

Short-Time Fourier Transform

The Short-Time Fourier transform (STFT) is a Fourier related transform used to determine the sinusoidal frequency and phase content of local sections of a signal as it changes over time [7]. Computing STFT requires the signal to be divided into segments of equal length and then have the FT applied to that segment [1]. This allows us to view the Fourier spectrum at a more granular level, which could potentially reveal different patterns amongst signals of different classes.

We can explain STFT through an example. Let $x(n)$ be a signal, and let $X_n(e^{j\omega_k})$ be the short-time Fourier transform of $x(n)$ at time n and frequency ω_k (Figure 4.2) [1].

$$X_n(e^{j\omega_k}) = \sum_{m=-\infty}^{\infty} w(n-m) x(m) e^{-j\omega_k m}.$$

Figure 4.1: Example of short-time fourier transform formula [1].

One of the issues with this method is that we can create very narrow window sizes, giving us a better understanding of the data with respect to time but losing an understanding of the frequency domain as a whole. Additionally, selecting an appropriate window size

for segmenting the signal can be an arduous task that would require fine-tuning as well as increasing the dimensionality of a dataset.

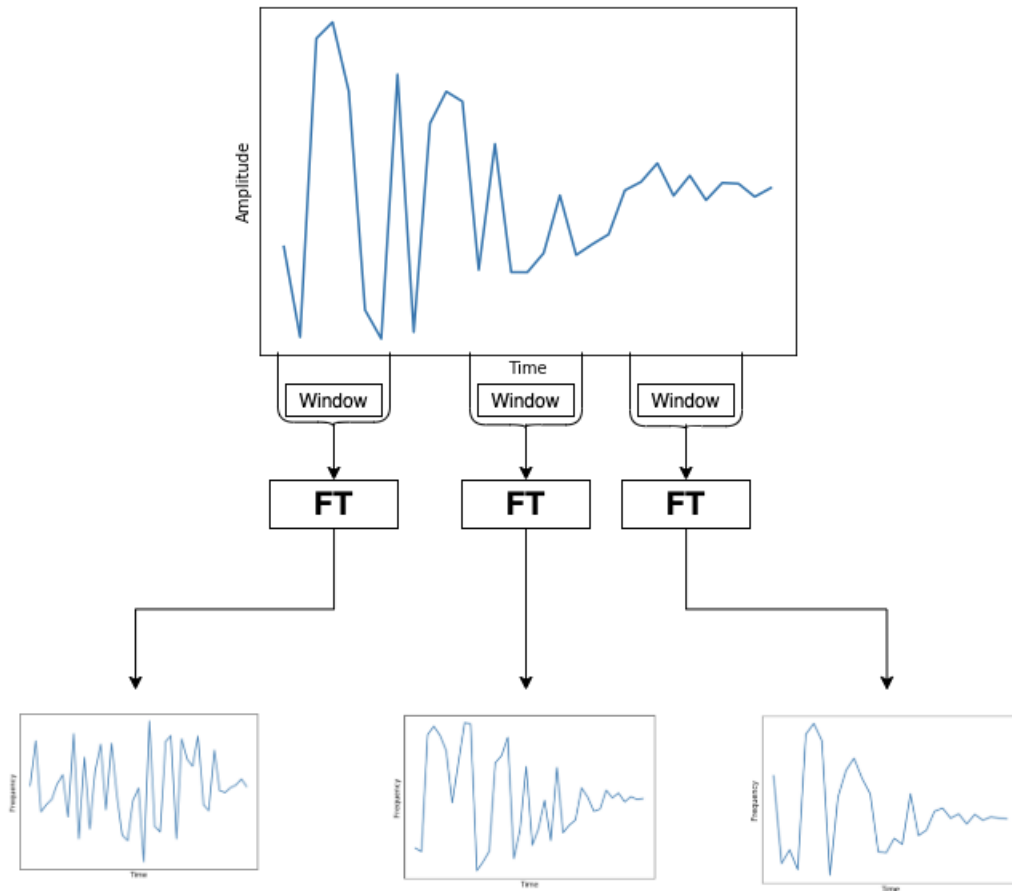


Figure 4.2: Example of Short-Time Fourier Transform.

Sliding Window Method

Similar to STFT, we segment our signal into N windows of equal length and perform FT on every segmentation. However, instead of keeping the output from FT, we only keep the mean value of each segmentation. This performed as well as STFT and gave us the advantage of reducing the dimensionality of our dataset.

Polar Coordinates

The representation of a complex number as a sum of a real and imaginary number, $z = x + iy$, is called its Cartesian representation. For every cartesian point, we calculate its radius and angular distance as real values and use those as features. This method outperformed the previous three methods but limited our feature size, which made it difficult to get consistent accuracy across all three of our datasets.

Wavelets

Wavelet transforms are similar to the FT in that they deconstruct a signal using representations of other signals [12]. The key difference is that wavelet signals are finite in time and frequency as opposed to sine and cosine signals, which can carry on indefinitely [15]. This allows us to extract information from a signal with respect to time and location.

A signal is analyzed using different versions of dilated and translated basis functions called the mother wavelet [15]. There are two types of wavelet transformations, discrete and continuous [8]. In this thesis, we focus on discrete wavelet transformations, which use a discrete set of wavelet scales and translations to decompose the signal into a mutually orthogonal set of wavelets.

We take advantage of wavelets by applying discrete wavelet transforms as a filter-bank, meaning they are composed of cascading high-pass and low-pass filters. This allows us to split a signal into several frequency sub-bands [15]. We focus on wavelet decomposition as our feature extractor because it outperformed all other methods in terms of accuracy and training time.

Data Shape

A neural network requires all data to be the same shape. However, because we downsample in wavelet decomposition, our data size is halved, making it impossible to place all levels of decomposition into the same dataset. In order to circumvent this, we tested two different methods: resampling the data and treating each level of decomposition as a unique dataset.

Resampling

We use spline interpolation to resize each level of decomposition into a single array because it allows to retain properties of our original data even in a higher dimension [6]. We tested different sizes of interpolation from N to $N/3$, where N is the size of the largest discrete signal size, and found no reduction in performance.

Each Level of Decomposition as a Dataset

Figure 4.3 depicts our model in which each level of decomposition is a unique dataset. This method and the resampling method both achieved similar results, so we chose to move forward with this method.

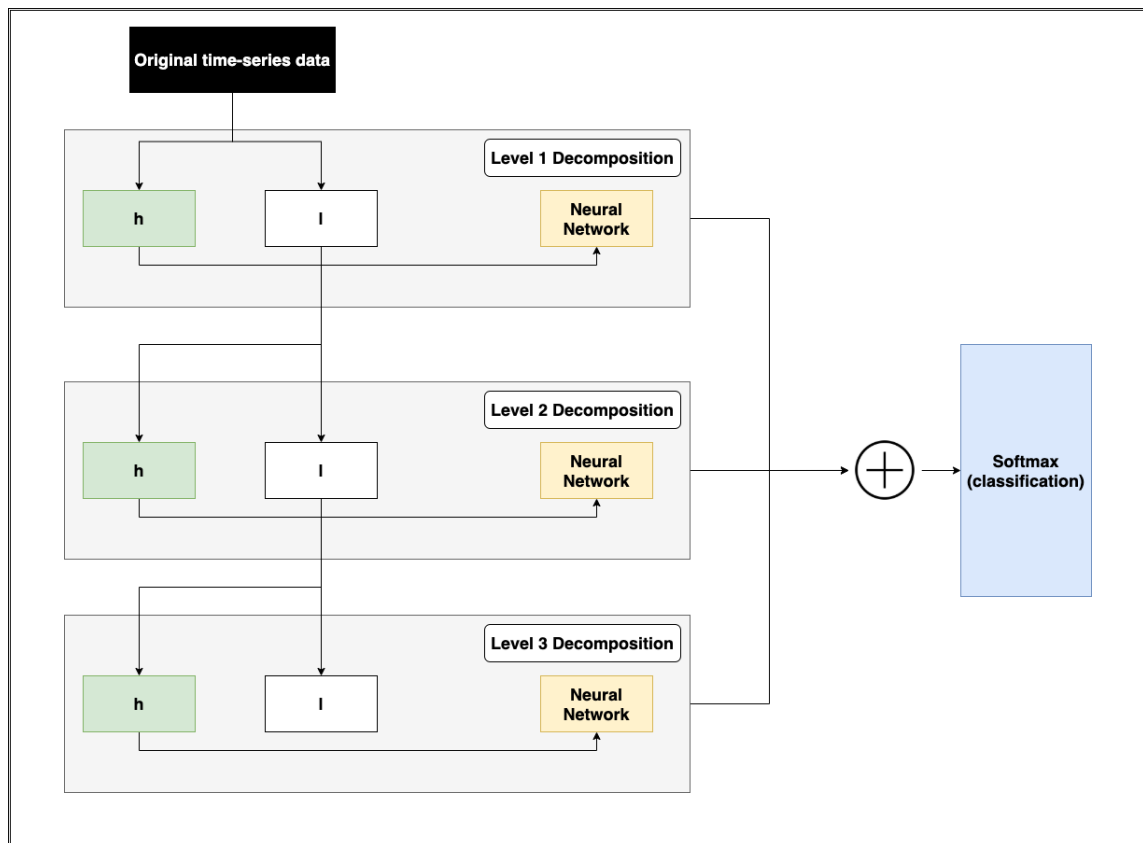


Figure 4.3: Example of a wavelet decomposition network.

Chapter 5

Models

Definition

Machine learning algorithms can be seen as learning a target function (f) that maps input data (X) to an output (Y). Then, f is used to make predictions on new input data. There are several techniques to make this work, but we focus on nonparametric algorithms, namely Convolutional Neural Networks (CNN) and Long Short-Term Memory networks (LSTM). Nonparametric algorithms attempt to make minimal assumptions about the form of the function to learn any form from data provided to it. An example would be a neural network as it has no prior knowledge of what it is classifying and attempts to generalize any new data points.

CNN

A Convolutional Neural Network is a type of deep neural network that can be applied to different domains such as computer vision or time series analysis. As seen in Figure 5.1, it consists of an input and output layer as well as several hidden layers. A convolution is an operation between a vector of weights w against an input x [19]. It consists of taking the dot product between m and x in steps of a filter size defined by n .

Convolutional Neural Networks perform feature learning via non-linear transformations implemented as a series of layers [3]. The input data is a multidimensional array, called a tensor. The tensor is passed through an input layer, followed by a series of hidden

layers to extract features, and an output layer, which in the case of classification, gives a probability for each class.

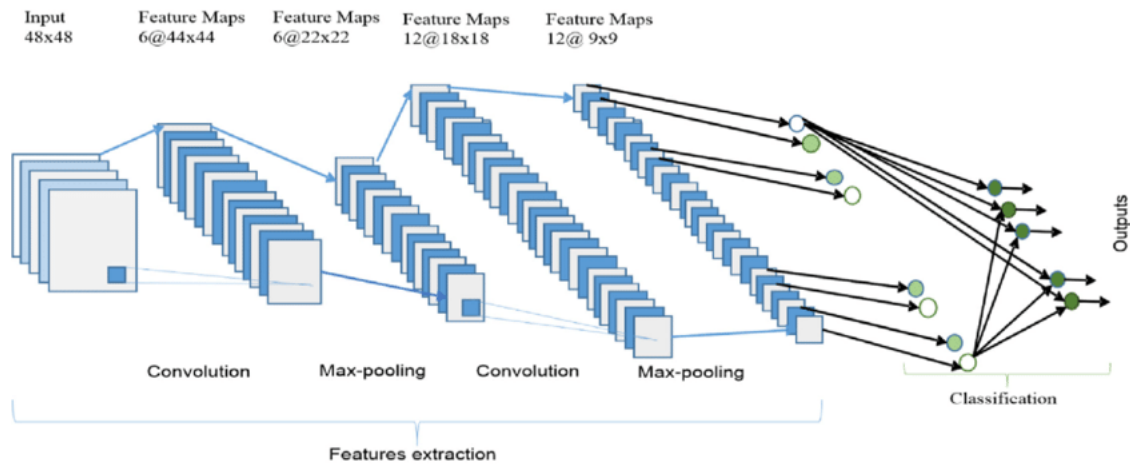


Figure 5.1: Example of a Convolutional Neural Network [2].

Hidden layers are crucial to neural networks because they help in determining which data representations are useful for explaining the relationships in the given data. Each layer consists of several kernels, which are feature detectors that convolve over the input data and output a transformed version of it.

LSTM

Long short-term memory (LSTM) is an artificial recurrent neural network (RNN) architecture used in the field of deep learning. Unlike standard feedforward neural networks, LSTM has feedback connections, processing not only single data points but also entire sequences of data.

A standard LSTM unit is composed of a cell, an input gate, an output gate, and a forget gate. The cell remembers values over arbitrary time intervals, and the three gates regulate the flow of information into and out of the cell [5].

LSTM networks are well-suited to classifying, processing, and making predictions based on time series data since there can be lags of unknown duration between significant events in a time series. LSTMs are developed to deal with the vanishing gradient problems that can be encountered when training traditional RNNs [5].

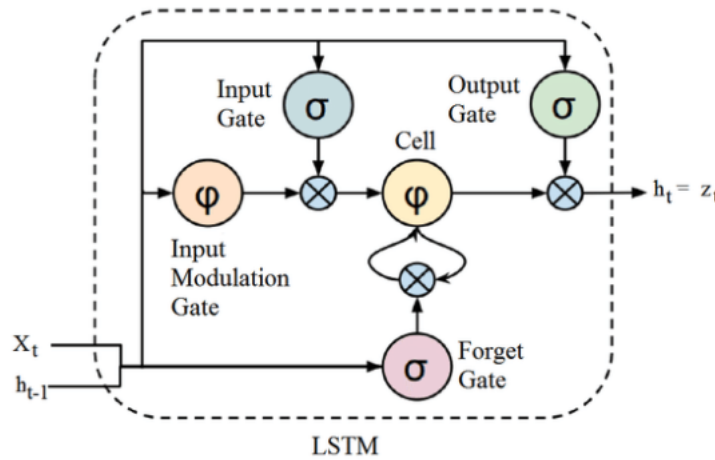


Figure 5.2: Example of a LSTM unit [5].

Multimodal Deep Learning

Due to the superior performance and computationally tractable representation capability in multiple domains such as visual, audio, and text, deep neural networks have gained tremendous popularity in multimodal learning tasks [9]. Typically, domain-specific neural networks are used on different modalities to generate their representations, and the individual representations are merged or aggregated [10]. Finally, the prediction is made on top of aggregated representation, usually with another neural network to capture the interactions between models and learn complex function mapping between input and output.

Architectures Used

CNN

Our CNN architecture consists of 3 convolutional layers each followed by a batch normalization and dropout layer then finally connected to a dense layer (Figure 5.3). We attempted to make it deeper but found inconsistent performance across all of our datasets.

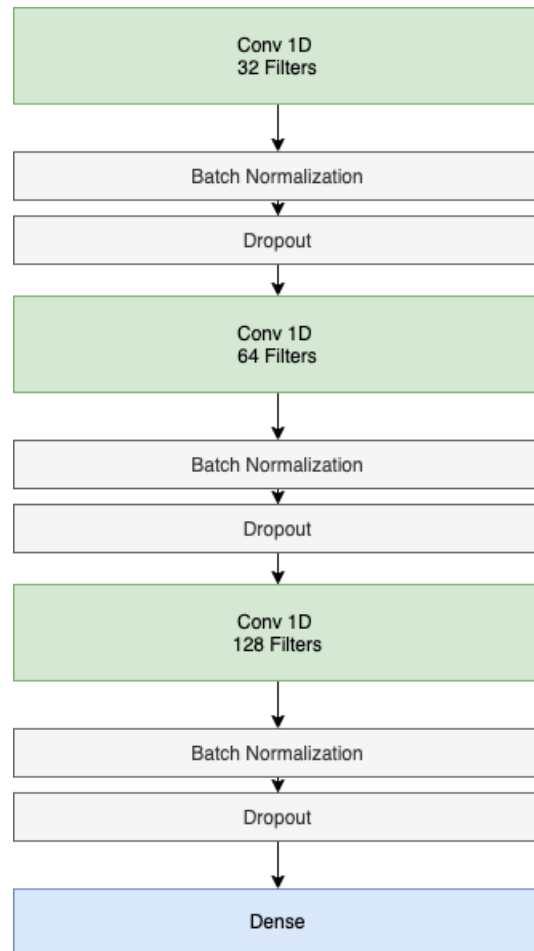


Figure 5.3: CNN architecture used in our black-box model.

LSTM

We wanted to keep our LSTM network as small as possible for simplicity, so we went with two LSTM layers of 128 hidden states and a dropout rate of 0.4 followed by a

dense layer of 128 connections. We attempted different state sizes but found 128 to be the smallest number of units with the highest and most consistent rate of accuracy.

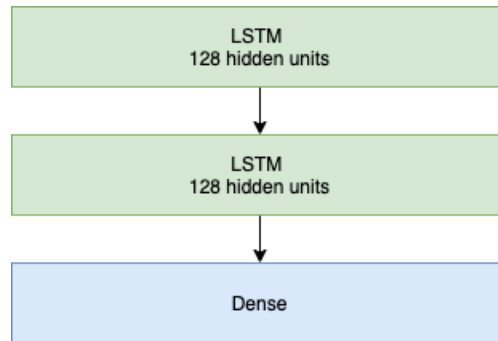


Figure 5.4: LSTM architecture used in our black-box model.

Final Model

We propose a multimodal learning architecture for time series classification, which is illustrated in Figure 5.5. The multimodal architecture allows our data to vary in size, and given that all three levels of decomposition have varying lengths due to downsampling, we propose that each model learn the representation of each distinct feature. For this model, we propose two architectures: a 1D CNN and an LSTM. We employ two CNN's and one LSTM, with each concatenated following their respective dense layer, after which classification is performed via a softmax layer.

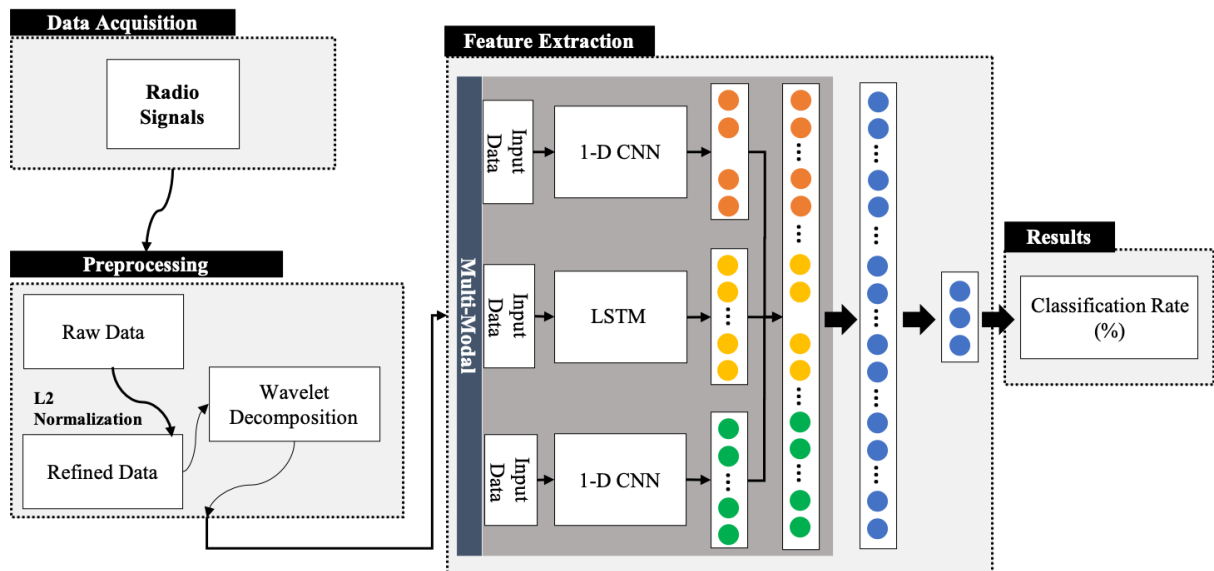


Figure 5.5: Overview of our black-box model.

Chapter 6

Markov Transition Fields

Method

We propose a framework similar to [18] for encoding dynamical transition statistics, but we continue their work by adding i^{th} order Markov transition probabilities.

Given a time series X , we decompose its magnitude axis into three separate properties, P_1, P_2, P_3 . We then identify the Q quantile bins for each property and assign each $x \in P_1, P_2, P_3$ to its corresponding bin q_j (j in $[1, Q]$). Thus, we construct three $Q \times Q$ adjacency transition matrices, W_1, W_2, W_3 , by counting transitions among quantile bins in the manner of an i^{th} order Markov chain along the time axis.

	A	B	C	D
A	0.917	0.083	0	0
B	0.083	0.583	0.334	0
C	0	0.260	0.522	0.218
D	0	0.083	0.167	0.75

Figure 6.1: Example of $Q \times Q$ quantile bin matrix used to calculate MTF [18].

W does not take into account the temporal axis, so to prevent any information loss, we construct a Markov Transition Field, M , for each W . The MTF denotes the probability of transitioning from q_i to q_j for each $x \in P$. This, in turn, allows us to consider the transition probability on the magnitude and temporal axis.

As described in [18] the MTF encodes the multi-span transition probabilities of the time series, but given that we have three different M 's, we modify their approach and consider each M to be a separate color channel of RGB where M_1 is red, M_2 is blue, M_3

$$M = \begin{bmatrix} W_{ij}|x_1 \in q_i, x_1 \in q_j & \dots & W_{ij}|x_1 \in q_i, x_n \in q_j \\ \vdots & \ddots & \\ W_{ij}|x_n \in q_i, x_1 \in q_j & \dots & W_{ij}|x_n \in q_i, x_n \in q_j \end{bmatrix}$$

Figure 6.2: Markov Transition Field where W_{ij} denotes the transition probability from quantile i to j .

is green. Since each row in M is a probability from 0 to 1, we multiply it by 255 to get a color value.

Classification

We apply MTF to our dataset and classify our images using a network similar to AlexNet, as shown in Figure 6.3. We avoided using the CNN used in our black-box model in order to create a deeper and higher performing network. For this model, we define our input as images with dimensions 224x224, use a kernel size of 3x3, and apply max-pooling after every convolutional layer. We have five convolutional layers using 11, 256, 256, 384, and 384 filters, respectively, followed by two dense layers and a softmax layer.

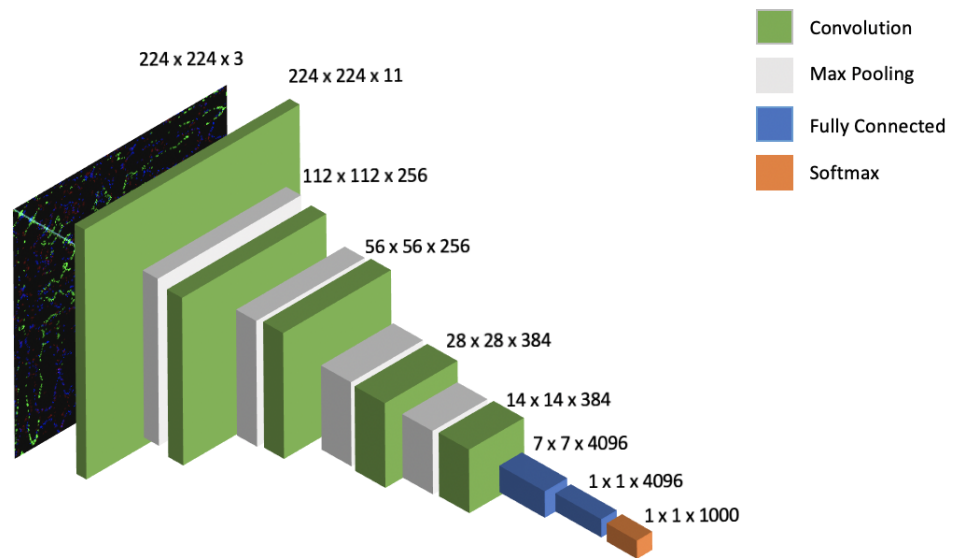


Figure 6.3: CNN architecture for classifying images generated through MTF.

Chapter 7

Main Results

Black-Box Model

Our black-box model is constructed of three different architectures: 1 LSTM and 2 CNN's that are all concatenated at their respective dense layer (Figure 5.5). The three features we use are the amplitude/phase of each level of decomposition for HV. We chose three because we wanted to increase the dimensionality of our data in the event that one had more descriptive features than the others.

Additionally, we tested our model using three levels of decomposition using the property VV as well as a combination of both VV and HV. We found no change in accuracy, therefore to limit the number of models we used, we only utilized the HV property.

To test our model, we split our data into training and test sets by randomly selecting 70% of our data as our training set and leaving the other 30% to be our test set. We then performed three-fold cross-validation and took the average of our results as the outcome. As can be seen in Figure 7.1, dataset A and C achieved similar accuracy, around 99%, and dataset B averaged at around 91% accuracy.

Dataset C

Given that this dataset is ten times larger than our other two datasets, we wanted to try different training and test splits, going as low as 8% for training data and using the remaining 92% as test data. We ran the same number of cross-validations as our previous tests and found minimal difference in performance.

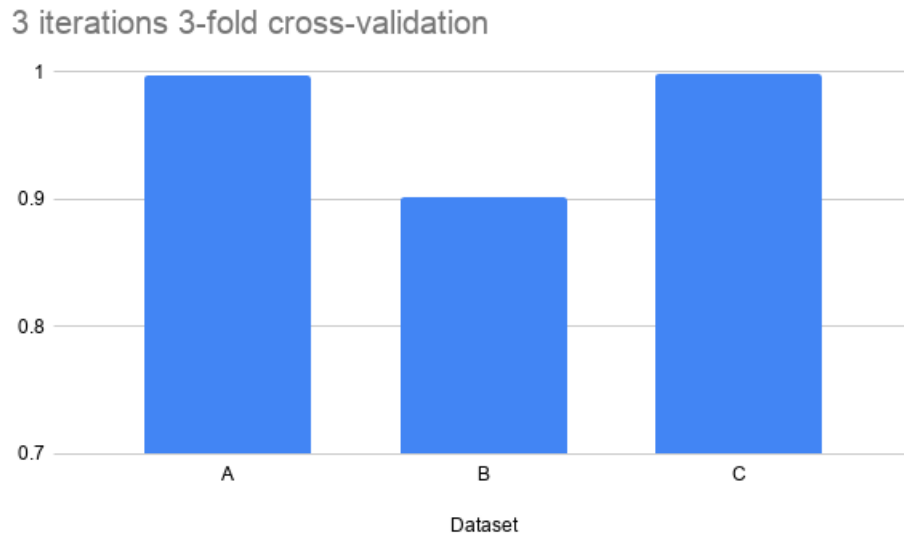


Figure 7.1: Accuracy on test data for our black-box model.

MTF

In order to generate images using MTF, we have to select three properties, and we chose P_1 as the radial value of HV's amplitude value, P_2 as the radial value of VV's amplitude value, and P_3 as the absolute squared distance between every P_1 and P_2 .

We then split those images using a randomized 70/30 split where 70% of the data is our training set, and the other 30% is our test set. Afterward, we trained our CNN (Figure 6.3) on those images and were able to achieve 60% accuracy across all three of our datasets and show the features amongst the classes behave differently, as can be seen in Figure 7.2.

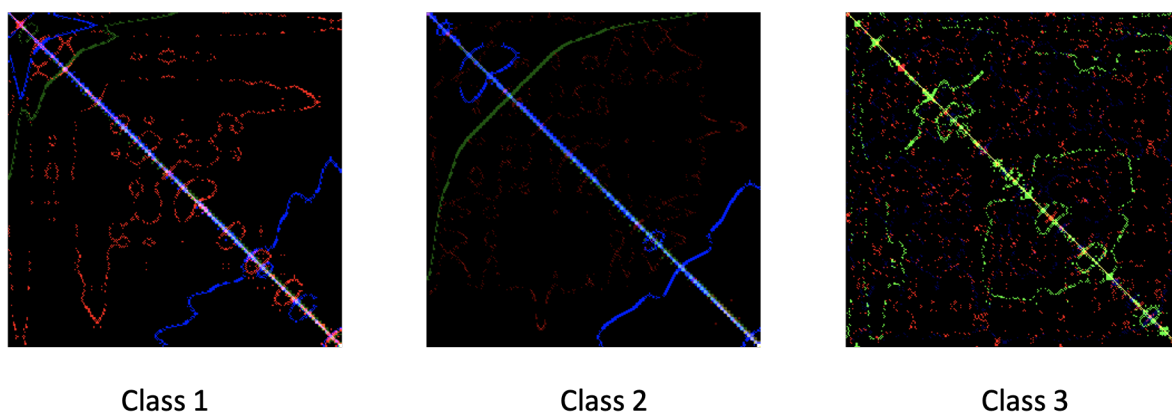


Figure 7.2: Examples of images generated through MTF.

Chapter 8

Conclusion

Throughout this thesis, we have stressed the importance of a hybrid machine learning algorithm as well as provided a framework that can successfully classify and visualize different classes when applied to time-series data.

We proposed an extension to existing work described in [18] to help users visualize their data while performing classification with 60% accuracy on that dataset and model. We also created a robust black-box model consisting of two different deep learning architectures that consistently provided competitive accuracy.

In the future, we would like to explore other methods for our white-box model that would act as an alternative to Markov Transition Fields in hopes that it produces higher accuracy than our current model. Additionally, we would like to extensively test our model on different time-series datasets outside of radio signal datasets.

Bibliography

- [1] J.B. Allen and L.R. Rabiner. A unified approach to short-time Fourier analysis and synthesis. *Proceedings of the IEEE*, 65(11):1558–1564, November 1977. Conference Name: Proceedings of the IEEE.
- [2] Md. Zahangir Alom, Tarek Taha, Chris Yakopcic, Stefan Westberg, Paheding Sidike, Mst Nasrin, Mahmudul Hasan, Brian Essen, Abdul Awwal, and Vijayan Asari. A State-of-the-Art Survey on Deep Learning Theory and Architectures. *Electronics*, 8:292, March 2019.
- [3] Muhammad Aqib, Rashid Mehmood, Ahmed Alzahrani, Iyad Katib, Aiiad Albeshri, and Saleh Altowaijri. Smarter Traffic Prediction Using Big Data, In-Memory Computing, Deep Learning and GPUs. *Sensors*, 19:2206, May 2019.
- [4] Zhicheng Cui, Wenlin Chen, and Yixin Chen. Multi-Scale Convolutional Neural Networks for Time Series Classification. *arXiv:1603.06995 [cs]*, May 2016. arXiv: 1603.06995.
- [5] Jeff Donahue, Lisa Anne Hendricks, Marcus Rohrbach, Subhashini Venugopalan, Sergio Guadarrama, Kate Saenko, and Trevor Darrell. Long-term Recurrent Convolutional Networks for Visual Recognition and Description. *arXiv:1411.4389 [cs]*, May 2016. arXiv: 1411.4389.
- [6] J. A. Gregory. Shape Preserving Spline Interpolation. June 1985.
- [7] Mitch Hill. THE UNCERTAINTY PRINCIPLE FOR FOURIER TRANSFORMS ON THE REAL LINE. page 17.
- [8] Maryam Imani and Hassan Ghassemian. Curve fitting, filter bank and wavelet feature fusion for classification of PCG signals. In *2016 24th Iranian Conference on Electrical Engineering (ICEE)*, pages 203–208, May 2016. ISSN: null.
- [9] Edward Kim and Kathleen F. McCoy. Multimodal Deep Learning using Images and Text for Information Graphic Classification. In *Proceedings of the 20th International ACM SIGACCESS Conference on Computers and Accessibility, ASSETS '18*, pages 143–148, Galway, Ireland, October 2018. Association for Computing Machinery.
- [10] Kuan Liu, Yanen Li, Ning Xu, and Prem Natarajan. Learn to Combine Modalities in Multimodal Deep Learning. *arXiv:1805.11730 [cs, stat]*, May 2018. arXiv: 1805.11730.

- [11] Raed Y. Mesleh, Harald Haas, Sinan Sinanovic, Chang Wook Ahn, and Sangboh Yun. Spatial Modulation. *IEEE Transactions on Vehicular Technology*, 57(4):2228–2241, July 2008. Conference Name: IEEE Transactions on Vehicular Technology.
- [12] Michel Misiti, Yves Misiti, Georges Oppenheim, and Jean-Michel Poggi. *Wavelets and their Applications*. John Wiley & Sons, March 2013. Google-Books-ID: Ee-MYyvA5PDoC.
- [13] Timothy J. O’Shea, Johnathan Corgan, and T. Charles Clancy. Convolutional Radio Modulation Recognition Networks. *arXiv:1602.04105 [cs]*, June 2016. arXiv: 1602.04105.
- [14] Julius Smith. Mathematics of the Discrete Fourier Transform (DFT): With Audio Applications - Julius Orion Smith - Google Books.
- [15] G. Strang and T. Nguyen. *Wavelets and Filter Banks*. Wellesley-Cambridge Press, 1996.
- [16] Jingyuan Wang, Ze Wang, Jianfeng Li, and Junjie Wu. Multilevel Wavelet Decomposition Network for Interpretable Time Series Analysis. *arXiv:1806.08946 [cs, eess, stat]*, June 2018. arXiv: 1806.08946.
- [17] Tong Wang and Qihang Lin. Hybrid Predictive Model: When an Interpretable Model Collaborates with a Black-box Model. *arXiv:1905.04241 [cs, stat]*, May 2019. arXiv: 1905.04241.
- [18] Zhiguang Wang and Tim Oates. Encoding Time Series as Images for Visual Inspection and Classification Using Tiled Convolutional Neural Networks. page 7.
- [19] Rikiya Yamashita, Mizuho Nishio, Richard Kinh Gian Do, and Kaori Togashi. Convolutional neural networks: an overview and application in radiology. *Insights into Imaging*, 9(4):611–629, August 2018. Number: 4 Publisher: SpringerOpen.