

Warsaw University of Technology

FACULTY OF  
MATHEMATICS AND INFORMATION SCIENCE



# Master's diploma thesis

in the field of study Computer Science  
and specialisation Data Science

Transfer learning for time series classification

**Paulina Pacyna**

student record book number 290600

thesis supervisor

Agnieszka Jastrzębska, PhD

WARSAW 2023



## Abstract

### Transfer learning for time series classification

The task of classifying time series is an important problem in the field of data mining. Most phenomena that change over time can be described in terms of time series. A time series can describe, for example, the amplitude of a heartbeat sound, stock prices, or hand movement along an axis when serving in tennis. Time series can express different characteristics of the phenomenon. Those characteristics are called *classes*. For example, the heartbeat sound amplitude can represent (belong to) one of two classes: *healthy* and *unhealthy*. Time series classification attempts to learn the distinctive features of a time series and build a model that can distinguish between those classes.

The time series classification problem was initially solved using classical algorithms such as the k-nearest neighbor classifier with distance measures suited for time series, like Dynamic Time Warping. The advantages of using deep learning algorithms in the context of time series classification have recently begun to be recognized. Neural networks can detect shapes that distinguish a class or understand ordered temporal relationships.

Transfer learning is practically used when there is limited data to train. Transfer learning attempts to apply patterns learned from one dataset to improve learning when creating a model for another dataset. A common practice is to prepare a source classifier trained on a large, readily available amount of data for one task and then use this model or parts of it for a detailed task with a smaller amount of data. Models trained using this method often have a shorter training time, faster accuracy increase, and can generalize more easily on the test set.

**Keywords:** time series, classification, transfer learning, deep learning, ...



## Streszczenie

Zastosowanie techniki transfer learning w zadaniu klasyfikacji szeregów czasowych

Zadanie klasyfikacji szeregów czasowych jest ważnym problemem w dziedzinie eksploatacji danych. Szeregi czasowe występują za każdym razem, gdy chcemy zmierzyć jakieś zjawisko, które zmienia się w czasie. Szereg czasowy może opisywać np. amplitudę dźwięku bicia serca, ceny akcji czy ruch ręki wzdłuż osi podczas serwisu w tenisie. Takie szeregi czasowe mogą wyrażać różne cechy zjawiska. Te cechy nazywane są *klasami*. Na przykład amplituda dźwięku bicia serca może reprezentować (należać do) jednej z dwóch klas: *norma* i *choroba*. Klasyfikacja szeregów czasowych polega na rozpoznawaniu charakterystycznych cech szeregu czasowego i budowaniu modelu, który potrafi rozróżniać klasy.

Problem klasyfikacji szeregów czasowych był początkowo rozwiązywany za pomocą klasycznych algorytmów, takich jak algorytm k-najbliższych sąsiadów w połączeniu z miarami podobieństwa dla szeregów, jak odległość DTW. W ostatnim czasie zaczęto dostrzegać zalety stosowania algorytmów głębokiego uczenia w kontekście klasyfikacji szeregów czasowych. Sieci neuronowe są w stanie wykryć kształty wyróżniające daną klasę lub zrozumieć relacje między obserwacjami w czasie.

Metoda transfer learning jest stosowana w przypadku ograniczonej ilości danych do trenowania modelu. Technika transfer learning próbuje zastosować wzorce wyuczone z jednego zbioru danych, podczas tworzenia modelu dla innego zbioru danych. Częstą praktyką jest przygotowanie źródłowego klasyfikatora wytrenowanego na dużej, łatwo dostępnej ilości danych dla jednego zadania, a następnie wykorzystanie modelu lub jego części do szczegółowego zadania z mniejszą ilością danych. Modele wytrenowane tą metodą często mają krótszy czas trenowania, szybszy wzrost dokładności i lepsze, ogólniejsze wyniki na zbiorze testowym.

**Słowa kluczowe:** szeregi czasowe, klasyfikacja, transfer learning ...



# Contents

<b>1. Introduction</b>	<b>11</b>
<b>2. Related works</b>	<b>13</b>
2.1. Time series classification	13
2.1.1. Multi Layer Perceptron	13
2.1.2. Convolutional Neural Networks	14
2.1.3. Fully Convolutional Networks	16
2.1.4. Residual Network	17
2.1.5. Encoder Network	18
2.2. The UCR archive	19
2.3. Preprocessing and augmenting time series	23
2.4. Transfer learning	23
2.4.1. Transfer learning categorization	24
2.4.2. Characteristics of a good source domain	24
2.4.3. Negative transfer in naive transfer learning approach	26
2.4.4. Inter-dataset similarity measure based on Dynamic Time Warping Barycenter Averaging	27
<b>3. Metodology</b>	<b>30</b>
3.1. Data reading	30
3.2. Preprocessing and preparing the data	30
3.3. Neural network architectures	31
3.4. Multi-source transfer learning	32
3.4.1. Baseline approach	32
3.5. Ensemble approach	34
3.6. Selecting the source datasets for multi-source transfer learning	36
<b>4. Results</b>	<b>38</b>
4.1. Experiments	38
4.1.1. Evaluation of the baseline approach	38





# 1. Introduction

Time series classification was initially approached using classical machine learning algorithms. Bagnall et al. in [1] categorizes the commonly used algorithms into several categories. The first category is time domain distance-based algorithms. Those algorithms use various distance measures adjusted for the time series domain to capture the similarity between pairs of time series. The distance measure is then combined with a distance-based classifier. A flagship example of an algorithm belonging to this class of algorithms is Dynamic Time Warping distance with the k-Nearest Neighbour classifier. The DTW-1-NN classifier is often used as a benchmark classifier. Another category of classifiers mentioned in [1] are dictionary, shapelet and interval based classifiers. All try to extract distinctive features for the time series class. The dictionary-based classifiers encode the time series into a dictionary of *words* representing the time series. Shapelet-based algorithms focus on subseries of the time series that are discriminatory of class membership. Interval-based algorithms extract features from selected intervals of the time series.

With the rise of the popularity of deep learning algorithms, researchers attempted to replace the former, hand-extracted features and algorithms with deep learning classifiers. Fawaz et al. reviewed deep learning algorithms applied to the time series classification task [7]. The usage of deep learning algorithms enabled the possibility to utilize transfer learning.

Transfer learning is widely used in image recognition and natural language processing. Fawaz et al. extended their previous findings by a study on the application of transfer learning [7]. The authors examine knowledge transferability on 85 datasets from the UCR archive by pre-training a model on one dataset and fine-tuning it on another. The authors examine if the accuracy improved for all pairs of datasets in the UCR archive.

Transfer learning for deep learning is usually done by training the network on a big, diverse dataset and utilizing the first layers of the network when training on another dataset. In image recognition or natural language processing, it is common to pre-train the source network on the ImageNet dataset or Wikipedia dataset, respectively. Such a diverse dataset with labels does not exist for time series. In this thesis, we will attempt to create an artificial, diverse dataset from 85 smaller datasets available on the UCR archive. We will try to mimic the good properties of a transfer learning source dataset by preprocessing, upsampling, and augmenting the dataset. We

will also study if transfer learning on a diverse dataset helps to generalize new training data and improve the training process on the target dataset. Similarly, as in [8], we will experiment to find which features of the source dataset (classes diversity, dataset size, augmentation, preprocessing) are fundamental for the transfer learning process.

## 2. Related works

In this chapter, we describe several algorithms used in time series classification, focusing on deep learning algorithms. We will refer to literature and papers in the domain of time series classification and transfer learning. We will also describe the most popular archive of time series datasets, the UCR archive.

### 2.1. Time series classification

A time series is an ordered collection of observations indexed by time.

$$X = (x_t)_{t \in T} = (x_1, \dots, x_T), \quad x_t \in \mathbb{R}$$

The time index  $T$  can represent any collection with the natural order. We assume that indices are spaced evenly in the set  $T$ . The realization or observation  $x_t$  in the times series is a numerical value describing the phenomena we observe, for example, the amplitude of a sound, stock price, or y-coordinate. Time series classification is a problem of finding the optimal mapping between a set of time series and corresponding classes.

#### 2.1.1. Multi Layer Perceptron

The Multi Layer Perceptron (MLP) is the first artificial neural network architecture proposed in [7] and can be used for time series classification tasks. Formally, the MLP network can be defined as a composition of functions (called *dense layers* or *layers*). The network returns a vector that usually represents the probability distribution over the set of classes.

$$MLP(X; \theta_1, \dots, \theta_M, \beta_1, \dots, \beta_M) = L_M(\dots L_2(L_1(X; \theta_1, \beta_1); \theta_2, \beta_2); \theta_M, \beta_M)$$

Each layer  $L_i : \mathbb{R}^M \rightarrow \mathbb{R}^N$  is a function that depends on the parameters  $\theta \in \mathbb{R}^{M \times N}, \beta \in \mathbb{R}^N$

$$L_i(X; \theta_i, \beta_i) = f_i(X\theta_i + \beta_i)$$

The function  $f_i : \mathbb{R}^N \rightarrow \mathbb{R}^N$  is an arbitrarily chosen non-linear function. The number of layers and dimension of weights in hidden layers is also arbitrary. The weights in the first and last

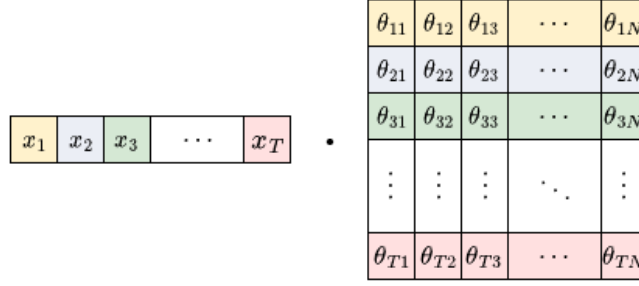


Figure 2.1: The multi layer perceptron architecture treats each entry in the input time series separately.

layer have to match the dimensionality of input data (e.g. the length of the time series) and the number of classes. The output of the last layer is interpreted as a probability distribution over the set of classes.

The disadvantage of using Multi Layer Perceptrons for time series classification is that the input size is fixed. All time series in the training data must have the same length. In transfer learning, this means that if we want to reuse the source network (or a set of first layers from the network), the time series in the target dataset must have the same length as in the source dataset.

The MLP architecture fails at understanding the temporal dependencies [7]. Each input value in the time series is treated separately because it is multiplied by a separate row in the weight matrix (see figure 2.1.1 for an illustration).

### 2.1.2. Convolutional Neural Networks

Convolutional Neural Networks are widely used in image recognition. A convolution applied for a time series can be interpreted as sliding a filter over the time series. A convolutional layer is a set of functions called convolutions or filters. The filter is applied at a given point, using the values surrounding the point.

To define the convolution operation, let's assume the input is a matrix  $X \in \mathbb{R}^{(N_1, \dots, N_K)}$ . In the case of images, the number of dimensions  $K$  is often equal to 3 (height, width, channels), for univariate time series we can assume just one dimension, and for multivariate time series we need two dimensions - (feature, time). The filter consists of a matrix of weights  $M \in \mathbb{R}^{(P_1, \dots, P_K)}$ . Usually,  $P_l$  are odd numbers, so that we can index the matrix with symmetrical numbers:  $(\frac{-P_l+1}{2}, \frac{-P_l+3}{2}, \dots, 0, \dots, \frac{P_l-1}{2})$ . The 0 index marks the center of the matrix.

## 2.1. TIME SERIES CLASSIFICATION

Finally the convolution  $*$  is defined as follows:

$$(X * M)_{i_1, \dots, i_K} = \sum_{l_1 = \frac{-P_1+1}{2}}^{\frac{P_1-1}{2}} \cdots \sum_{l_K = \frac{-P_K+1}{2}}^{\frac{P_K-1}{2}} M_{l_1, \dots, l_K} X_{i_1+l_1, \dots, i_K+l_K}$$

The result of the convolution is passed elementwise to a nonlinear function. The nonlinear function, together with the convolution operation, will be called a filter. In the case of univariate time series, the first layer of the convolutional neural network is one-dimensional. The output of the first layer has dimensions (length of time series - the length of the filter + 1, number of filters). Below we define the value of the output for filter  $i$

$$y_{t,i} = f_i([\theta_{\frac{-M+1}{2}}^i, \dots, \theta_{\frac{M-1}{2}}^i] \cdot [X_{t+\frac{-M+1}{2}}, \dots, X_{t+\frac{M-1}{2}}]),$$

where  $\cdot$  is a dot product and  $t \in T$  is the time index. In figures 2.2 and 2.3 we show the results of applying two different filters to a time series.

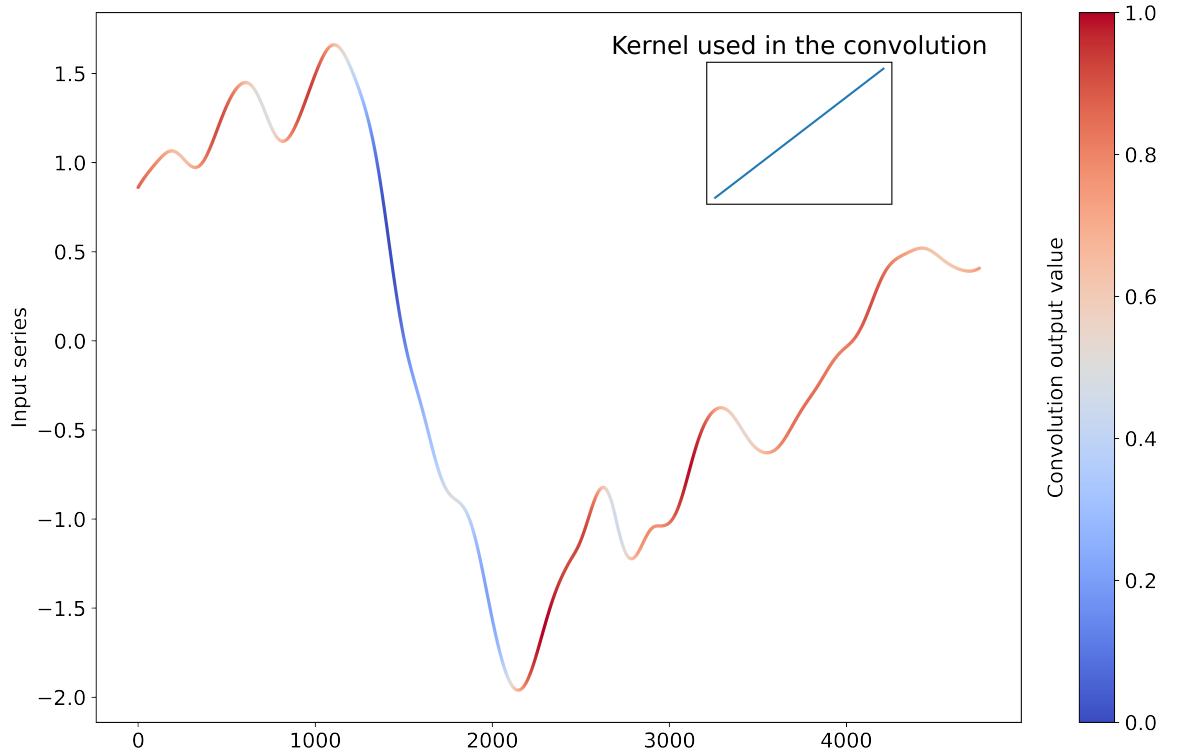


Figure 2.2: Result of applying a "slope" detecting filter to time series. Red regions indicate a higher value in the time series returned by the output convolution.

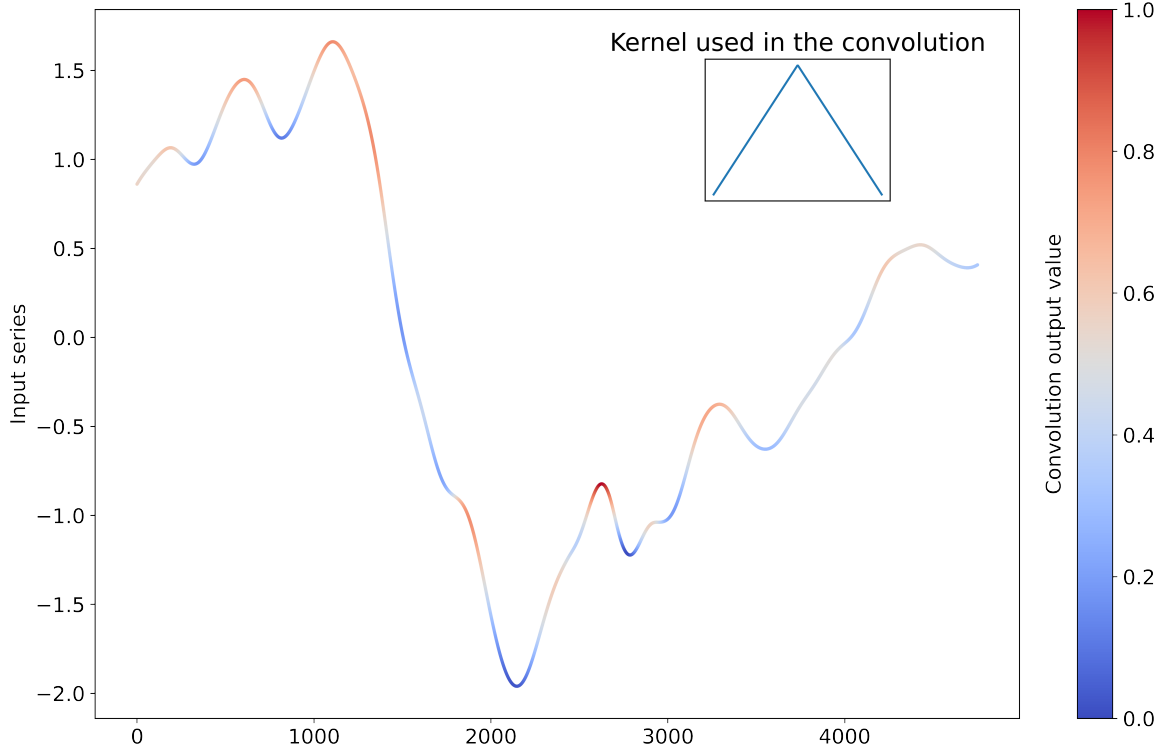


Figure 2.3: Result of applying a "peak" detecting filter to a time series. Red regions indicate a higher value in the time series returned by the output convolution.

The weights  $\theta^i$  are different for each filter. The same filter is applied over the whole length of the time series. This is called *weight sharing*, enabling the network to learn patterns regardless of the position in the time series.

The architecture of the convolutional layer is not dependent on the input data size. Regardless of the size of the input data, the number of filters and size of filters remain the same, and only the output sizes depend on the input size. Therefore, if the convolutional layer is succeeded by layers with the same property, like other convolutional layers or Global Pooling with Dense Layer (see Section 2.1.3), the whole network may be invariant to the input sizes [7]. Such networks may be interesting in transfer learning, as the sizes of time series in the source and target tasks do not have to match.

### 2.1.3. Fully Convolutional Networks

Fully Convolutional Networks are convolutional networks used in time series classification. A sample architecture of a Fully Convolutional Network proposed is in [7]. The first layers in the network are 3 blocks of convolutional layers with ReLU activation function followed by batch normalization layers. The output of the last block is passed to a global pooling layer. The global

## 2.1. TIME SERIES CLASSIFICATION

pooling layer averages the output through the time axis, resulting in a vector of length equal to the number of feature maps in the last convolutional layer. The averaged vector is passed to a block of 2 fully connected layers. Figure 2.4 shows a visualization of the network.

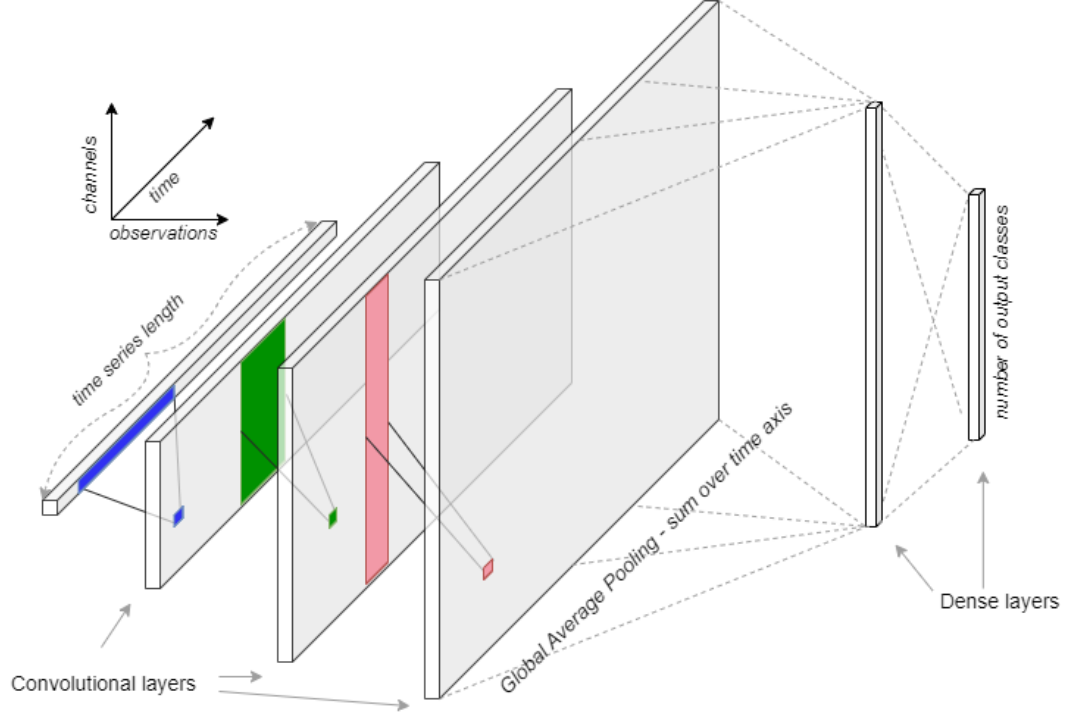


Figure 2.4: Architecture of a Fully Convolutional Network. Source: [7]

Because the architecture of convolutional layers does not depend on the size of input data and the convolutional layers are followed by pooling over the time axis, the whole network is capable of processing data of variable lengths.

### 2.1.4. Residual Network

The Residual Network was first proposed in [13]. The Residual Network is a relatively deep architecture compared to other neural networks used for time series classification.

A residual connection addresses the vanishing gradient problem occurring in networks composed of many layers [13]. The vanishing gradient problem occurs in the backpropagation process. When the gradient passes from the last layer to the first, it may decrease toward zero. This causes the first layers of the network to learn slowly. The residual connection between layers passes the input directly from one layer to another, skipping a few layers. This way, if the network is struggling with a vanishing gradient, this shortcut connection may help the network to converge.

The architecture of the Residual Network is conceptually similar to the FCN network. Instead of three convolutional layers, the Residual Network begins with three blocks of convolutional

layers, connected with residual connections. Each block consists of three convolutional layers. The three blocks, which consist of nine convolutional layers and three residual connections in total, are followed by a Global Average Pooling layer and a Dense layer. The networks' diagram is shown below (Figure 2.5).

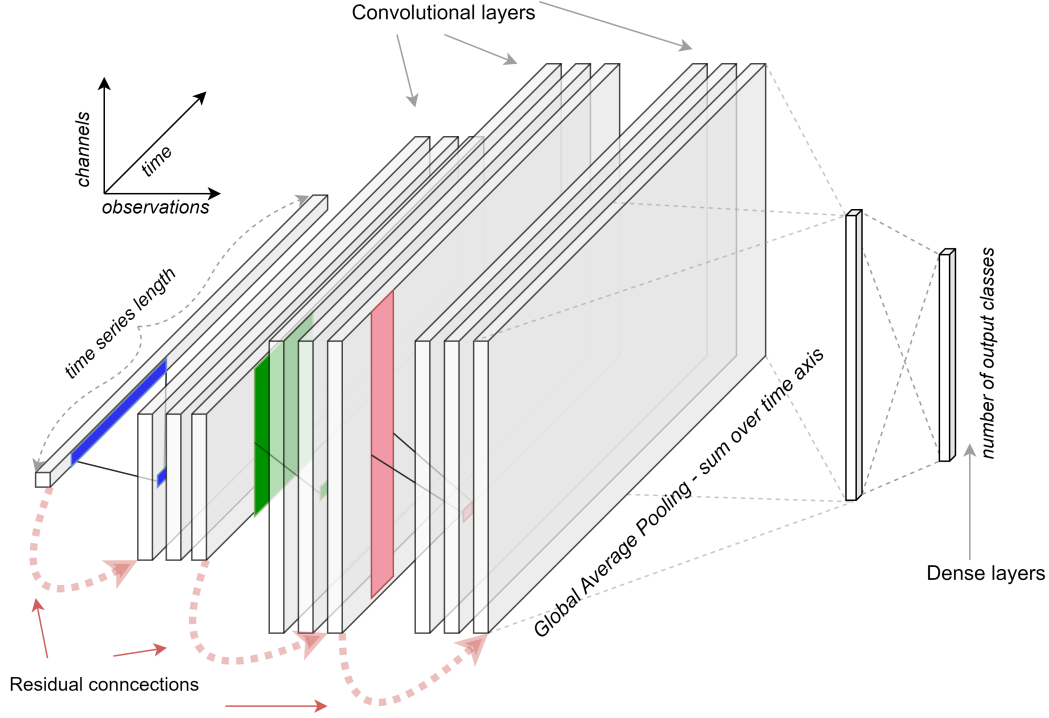


Figure 2.5: Architecture of a Residual Network. Source: [7]

### 2.1.5. Encoder Network

The Encoder Network is a deep convolutional network proposed by [11]. The first layers of the network are similar to the FCN architecture. The network consists of three convolutional layers followed by an attention layer instead of the Global Average Pooling layer. Another novelty introduced in this network compared with the FCN architecture is adding a Dropout layer and replacing the ReLU activation function with PReLU (Parametrized ReLU), thus adding another parameter to the network. The network ends with a Dense layer predicting the distribution over the classes. The architecture is shown below (Figure 2.6).



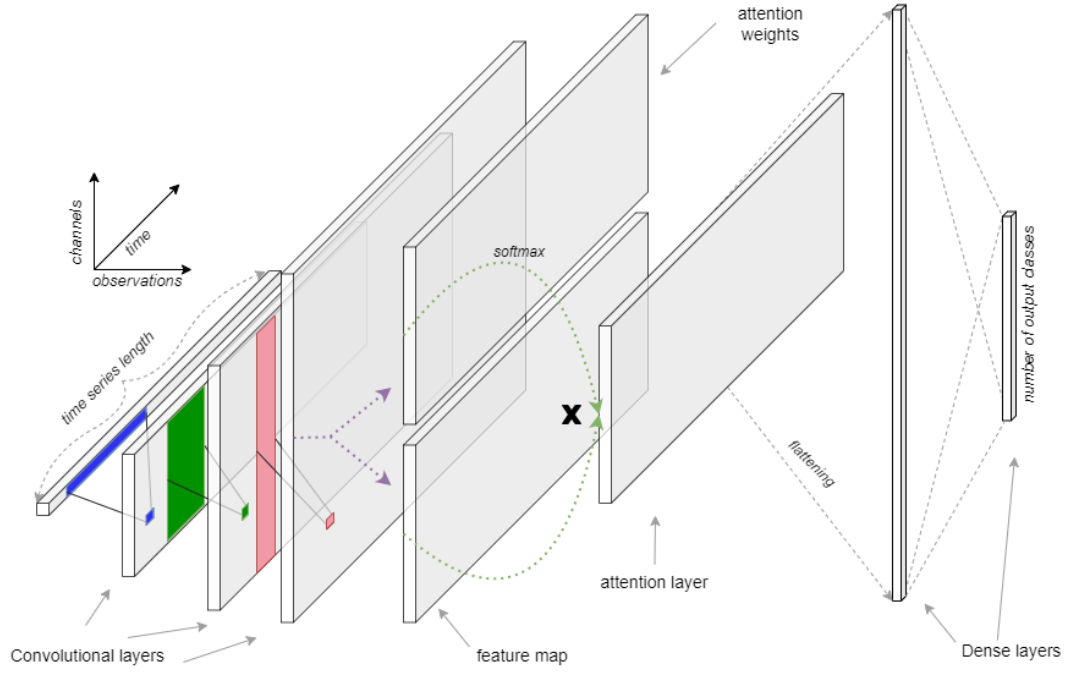


Figure 2.6: Architecture of an Encoder Network. Source: [7]

The attention layer, first proposed in [2], assigns weights to each input element representing the importance of the prediction. The weights representing the attention are normalized using the softmax function and then applied to the input data. The weights are learned by the network. The outputs of the attention layer can be interpreted as a universal representation of the input time series that can adapt to and represent unseen data [11].

As opposed to the FCN network architecture, the Encoder network is not invariant to the input size. The Dense layer parameters depend on the output size of the attention layer and implicitly on the output size of convolutional layers. Similarly, as in the former architecture, the convolutional layers enable weight sharing, thus learning shapelets independently of their position in the time series.

## 2.2. The UCR archive

The UCR (University of California, Riverside) archive is a time series classification archive introduced in 2002 by Dau, Keogh et al. [5]. At moment of writing this thesis, it consists of 140 univariate time series datasets. The time series describes various phenomena, like readings from a device or sensor, motion recordings, or food spectrographs. The summary of the time series types is listed below (Table 2.1). We also show a few time series samples from different datasets below (Figure 2.7).

Dataset type (source)	Number of datasets
Audio	10
Device	10
ECG (Electrocardiogram)	7
EOG (Electrooculography)	2
EPG (Electrical penetration graph)	2
Hemodynamics	3
Image	33
Motion	28
Sensor	22
Simulated	9
Spectro	12
Traffic	2

Table 2.1: Summary of types of datasets in the UCR archive.

All datasets together contain 1122 unique labels. Most of the datasets differentiate between two labels and have several dozen rows per class. The series length varies from 8 to almost 3 thousand. Below we show four histograms summarizing the number of classes per dataset, the lengths of the series, and the number of observations per class and per dataset. (Figure 2.8) The datasets in the archive are already split into train and test datasets. Most of the time series are already z-normalized. However, the authors recommend testing on additional test/train splits and preprocessing the time series accordingly to the method used. We will now proceed to describing the datasets in detail, considering all available categories.

- **Audio** - Sound amplitudes of heartbeats and spectrograms of sound produced by insects, cats, dogs or whales. Classes represent normal and abnormal patients or different species.
- **Device** - This category includes time series describing power consumption of house devices, or summarized population consumption for certain periods. This category is characterized by the a combinations of spikes and constant areas in the time series. The data is divided to classes based on different types of devices or different time periods of data collection.
- **ECG** - The data is composed of electrocardiogram readings. The classes corresponds to different persons being examined, different diseases detected, location of sensors or time when the data was collected.
- **EOG** - There are two datasets describing electrooculograph readings, measuring the corneo-retinal standing potential that exists between the front and the back of the hu-

## 2.2. THE UCR ARCHIVE

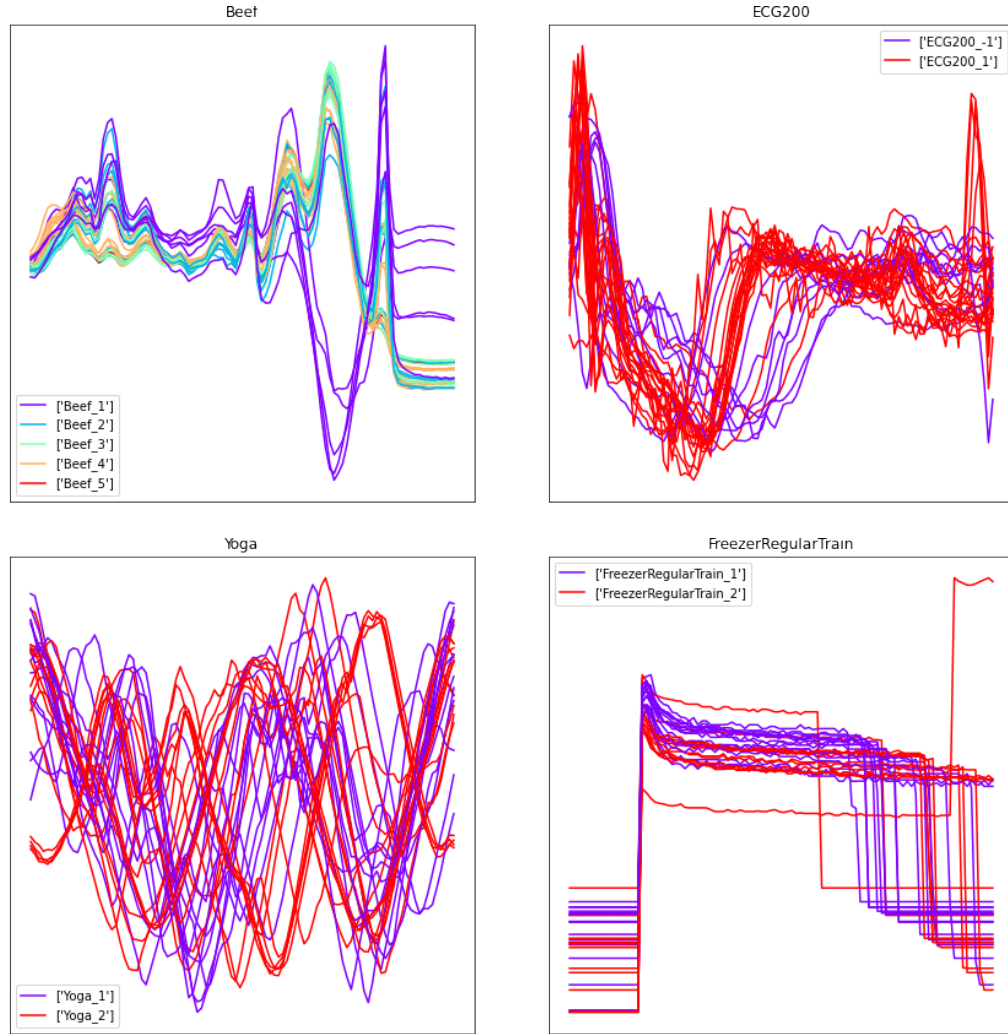


Figure 2.7: Example time series from the UCR archive. The images belong to the following categories: Food, ECG, Image, Device.

man eye. The classes distinguish how the reading varied when writing and looking at different characters.

- **EPG** - The electrical penetration graphs measure voltage changes of the electrical circuit that connects insects and their food source.
- **Hemodynamics** - This category consists of time series representing airway, arterial blood, central venous and pressure measurements. Each class represent one of the pigs that were examined.
- **Image** - Majority of datasets in this category are outlines of images of different shapes. The images represents algae, arrows, birds, bones, faces, vegetables, written words, fish, hands and leafs. Another types of time series present in this category are changes of pixel

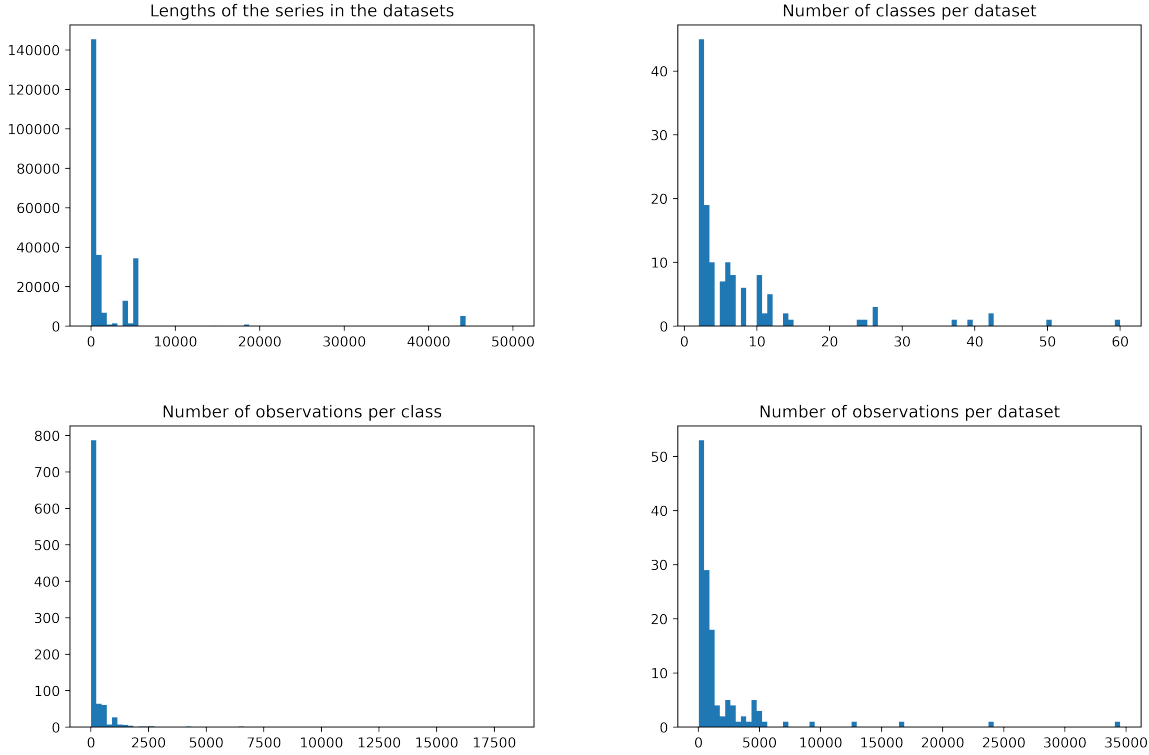


Figure 2.8: Summary of the UCR archive. The first plot is a histogram of time series length. The second histogram shows the number of classes in datasets. The two last histograms summarize the number of rows per class and per dataset.

intensity over time, or growth of an area on a picture over time.

- **Motion** - Positions of a certain points (usually on human body, but also on a car or worm) over an axis, e.g while walking, speaking, driving, playing cricket, skating, writing or performing gestures. One of the datasets uses electromiographic signal instead of point positions. Another dataset is different from others, as it describe number of pedestrians on different hours.
- **Sensor** - The data includes readings from sensors, e.g. earthquake sensors, humidity, temperature sensors or robots. An example task is to predict an earthquake or to classify the surface the robot was walking on.
- **Simulated** - The time series used in this category were simulated or obtained from an algorithm.
- **Spectro** - This category consist mostly of food spectrograms but also one data set of time series representing surface electromyography which captures electric activity of groups of muscle at rest or during activity

- **Traffic** - In this category are two datasets recording pedestrian flow. The classes correspond to weekdays or places.

### 2.3. Preprocessing and augmenting time series

Data preprocessing is the process of preparing, cleansing, and manipulating the data to optimize the training process. Data augmentation is a process used to increase the amount and variability of data. Data augmentation is beneficial when there is a limited amount of data, as it can prevent overfitting.

Many classical and deep learning models were invented and designed together with the data preprocessing and augmenting step [1, 7]. One such method is Window Slicing [7]. This method extracts subsequences of the original series. The resulting subsequences are then concatenated with a downsampled copy and smoothed copy. This method is used in Multi-scale Convolutional Neural Network [3] and Time Le-Net [9].

Another preprocessing and augmentation method proposed together with the former method in the Time Le-Net model was Window Warping. It aims at making the model more robust to perturbations in the time axis by training the network on *squeezed* or *dilated* series together with the original series. The *dilated/squeezed* time series is two times longer/shorter than the original time series. The *dilated*, *squeezed*, and original time series are concatenated, forming a vector 3.5 longer than the original time series. In the Time Le-Net model, this vector undergoes the Window Slicing preprocessing step before training.

### 2.4. Transfer learning

Transfer learning is a technique that attempts to apply knowledge learned while solving one task to enhance the learning process for another task. Formally, the problem can be described using the notions of tasks and domains [14, 16]. A *Domain* is a pair  $\mathcal{D} = (\mathcal{X}, P(\mathcal{X}))$ , where  $\mathcal{X}$  is the feature space (e.g., the time series observations, and  $P(\mathcal{X})$  is the probability distribution over the feature space. A *Task* is a pair of label space  $\mathcal{Y}$  and the decision function  $f$ ,  $\mathcal{T} = (\mathcal{Y}, f)$ . The decision function  $f$  is learned from  $\mathcal{X}, P(\mathcal{X}), \mathcal{Y}$  in the learning process.

Transfer learning attempts to utilize knowledge from different domain/domains and task/tasks. Formally, given  $S \in \mathbb{N}$  source domains and source tasks  $(\{(\mathcal{D}_i^S, \mathcal{T}_i^S) : i = 1, \dots, S\})$  and  $T \in \mathbb{N}$  target domains and target tasks  $(\{(\mathcal{D}_i^T, \mathcal{T}_i^T) : i = 1, \dots, S\})$  transfer learning utilizes

knowledge learned from source domains and tasks to improve the learning process of decision functions in target tasks  $\mathcal{T}_i^T$

#### 2.4.1. Transfer learning categorization

Transfer learning can be categorized from different points of view [16]. One of the ways in which we can divide transfer learning algorithms is based on the availability of the labels:

- **inductive** transfer learning - labels are available for both source and target datasets
- **transductive** transfer learning - labels are available only in the domain dataset
- **unsupervised** transfer learning - no labels available in either dataset

Another way of dividing transfer learning methods is by comparing feature space distribution and labels. If the source and training dataset consists of the same features, belonging to a similar distribution ( $\mathcal{D}^S = (\mathcal{X}, \mathcal{P}(\mathcal{X}))^S = (\mathcal{X}, \mathcal{P}(\mathcal{X}))^T = \mathcal{D}^T$ ) and the label spaces are equal ( $\mathcal{Y}^S = \mathcal{Y}^T$ ), then the task can be described as homogeneous transfer learning. If at least one of the former assumptions does not hold, the transfer learning setting is heterogeneous.

Transfer learning can be categorized based on the algorithm/approach used. There are four main categories [12] used in deep learning:

- **instance** based transfer learning - target dataset is enriched by instances from the source dataset belonging to the target distribution.
- **mapping** based transfer learning - source and target dataset are mapped to one domain. The distribution is adjusted in both datasets.
- **network** based transfer learning - target classifier uses parts of the network from the source classifier  $f$ .
- **adversarial** based transfer learning - an adversarial network tries to distinguish between the two datasets. If the network fails at this task, then it means that the extracted features are similar between source and target datasets.

#### 2.4.2. Characteristics of a good source domain

In the field of image processing, it is widespread to use convolutional neural networks pre-trained with the ImageNet dataset [8]. ImageNet is a large dataset of human-annotated images. It contains 1 million labeled images of 1000 classes. The label space consists of fine-grained classes such as breeds of dogs and cats and coarse-grained classes like *red wine* and *traffic*

*light*. Example pictures from this dataset are shown below 2.9. As transfer learning based on this dataset became more popular and successful, a question arose: Which features of this dataset make it so suitable for this task?



Figure 2.9: Sample images from the ImageNet dataset, with examples of similar (fine-grained) classes like two birds of different breeds and coarse-grained classes, like a lamp and a bird. Source: [url](#)

A study conducted in [8] attempts to answer this question. The first hypothesis is that the volume of the dataset is relevant to train accurate, general classifiers. The authors compared models that were pretrained on the original dataset and those based on sampled subsets (reduced 2, 4 8 and 20 times). The results show that the more training examples, the better results. The accuracy of the initial classifier occurred to be more dependent on the size of the dataset than the accuracy of classifiers fine-tuned from the base classifier. It is natural that the base classifier's accuracy will depend on the dataset size. Still, the classifier fine-tuned from the base classifier seems to cope with the reduced dataset, and the impact on the accuracy is less detrimental. This is visible in figure 2.10.

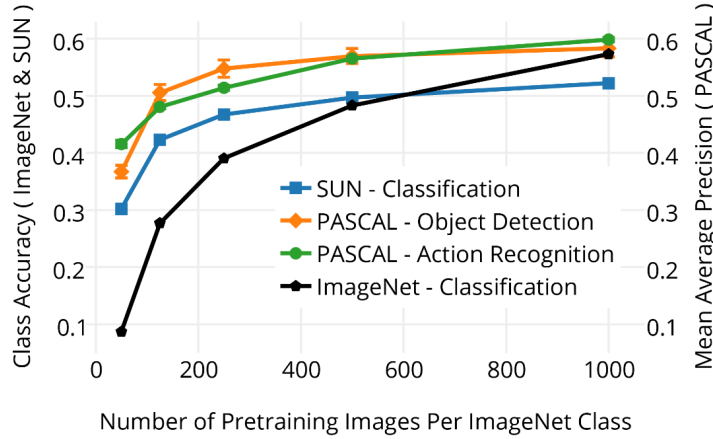


Figure 2.10: Accuracy of the base classifier (black) and classifiers fine-tuned from the base classifier. Image source: [8]

Next experiments examine the label space. The authors study if the granularity of the label space is essential for the problem. To compare the results, the label space is clustered, and 127 classes are derived from the initial 1000 classes. Pre-training with the reduced label space has a minimal negative impact on the accuracy of classifiers fine-tuned from this classifier. This suggests that such a fine division may not be needed.

Finally, the last question is if we train the classifier on the reduced label space with 127 classes, will it be able to distinguish between the fine-grained classes? To examine that, the authors extracted features from the first layers of the networks trained on reduced label space. Then, the authors performed classification with 1-NN and 5-NN models on the extracted feature space but with 1000 classes. The findings are that the k-NN classifier performs 15% worse on a reduced dataset versus the normal dataset.

While the article [8] does not conclude which single feature of the ImageNet dataset makes it so efficient as a source dataset for transfer learning, it is clear that all properties of this dataset are essential for the accuracy of the classifier. We will try to recreate the properties of the ImageNet dataset when creating the source dataset from time series.

### 2.4.3. Negative transfer in naive transfer learning approach

Negative transfer is a problem in the transfer learning process when the target classifiers achieve worse accuracy than without using transfer learning. Transfer learning for time series has been investigated in [6]. The paper's authors tested transfer learning on all 85 datasets from the UCR archive. In the first approach, called *naive transfer learning*, the authors used all pairs of datasets. One dataset from the pair was used as a source dataset to prepare the source



classifier for transfer learning. The second dataset was used to fine-tune the former classifier. The architecture used in this approach was the Fully Connected Network, described in section 2.1.3. All layers except for the last one were transferred from the source classifier to the target classifier. The authors compared the accuracy of a fine-tuned classifier to that of a classifier trained without transfer learning.

With respect to different source datasets, target classifiers' accuracies exhibited a high variance. There were cases where the fine-tuned classifiers suffered from transfer learning. On the other hand, for each target dataset, there was always a source dataset beneficial for the classification results. It is shown in the picture 2.11, where for each target dataset, the maximum, median, and minimum accuracies for the different experiments were plotted.

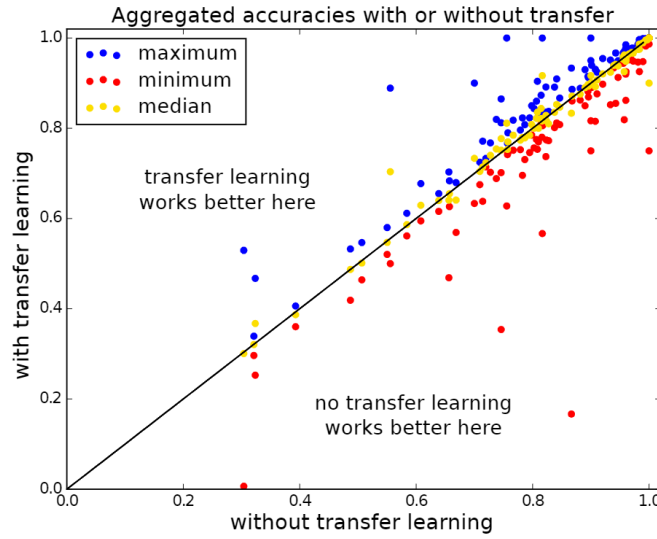


Figure 2.11: Image source: [6]

A negative transfer learning effect can be minimized or avoided by optimally choosing the source dataset/datasets. To improve the transfer learning process, the authors proposed a method of choosing the source dataset based on the DTW distance. This method will be described in 2.4.4.

#### 2.4.4. Inter-dataset similarity measure based on Dynamic Time Warping Barycenter Averaging

Choosing the source dataset for transfer learning randomly, or by trial and error, can lead to the *Negative transfer* problem, described in the previous chapter (2.4.3). In [6], the authors proposed an inter-dataset similarity measure, that is used to optimally choose the source dataset for transfer learning. The first step of the method is to compute a *prototype* time series per each class, for every dataset that is considered. This reduces the volume of the data, that will be

used is the next step. The next step is to compute the distances between each pair of datasets, choosing the minimum DTW distance between prototypes of each class in both datasets. We proceed to describe the algorithm in detail below.

The algorithm uses DTW barycenter averaging [10]. The result of this method is a times series that minimizes the sum of DTW distances to each time series from the dataset.

$$DTW\_average = \arg \min_{X \in \mathcal{X}} \sum_{\bar{X} \in \mathcal{X}} DTW(X, \bar{X})$$

The resulting time series does not need to belong to the original dataset  $\mathcal{X}$ .

The first step of the algorithm is the data reduction step. For each class in the dataset  $\{(X_i, y_i) : i = 1, \dots, N\}$ ,  $y_i = c \in \{1, \dots, C\}$  the prototype time series  $P_c$  is defined as the result of DTW barycenter averaging performed on the subset of  $\mathcal{X}$  corresponding to class  $c$ . As a result, the original dataset is reduced to  $C$  observations.

The next step is to compute the distance between the target dataset and all possible source datasets. The distance between two datasets is the minimum DTW distance between each pair of prototypes from the reduced datasets. The whole algorithm is described in detail in [6].

In [6] the distance was computed pairwise for all available datasets (using train splits). The authors compare the transfer learning performance with a randomly chosen source dataset versus choosing the source dataset that minimizes the former distance to the target dataset. The results are shown below (Figure 2.12).

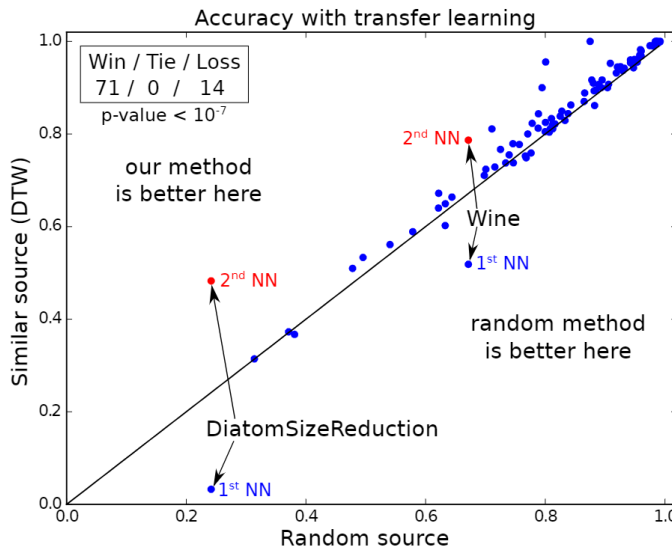


Figure 2.12: Image source: [6]

The result shows that choosing the source dataset for transfer learning based on similarity to the target dataset may lead to benefits. In this work we will try to construct the dataset

artificially, by combining, preprocessing, and augmenting the datasets from the UCR archive, to mimic the properties of the ImageNet dataset. The artificial dataset will be composed of multiple datasets. Therefore we rather won't use this similarity to choose one particular dataset. We will take the inter-dataset distance into account when adjusting the distribution of classes in the artificial dataset, to favour classes similar to the target dataset.

### 3. Metodology

#### 3.1. Data reading

The data downloaded from the UCR archive is available in two formats, `ts` (as text files) and `arff/weka` (as binaries). In order to speed up the reading and conversion time, all datasets were read, concatenated and saved as one array of time series (the original lengths were preserved). At this point, only basic preprocessing was done. It included the following steps:

- replacing NaN values with zeros
- removing zeros at the end of the time series. This step was done, because some of time series published in the UCR archive were already padded with zeros to have the same lengths. This assumption is not necessary when using Fully Connected Network.
- trimming the series if the length of a time series exceeded 50000 time steps.

All time series were concatenated to one array of series and saved in a `.npy` format.

#### 3.2. Preprocessing and preparing the data

In order to avoid focusing on one particular form of preprocessing and preparing the data, the preprocessing was done on the fly when training the data. This way, fast experimenting with hyperparameters and methods was possible. In order to achieve it, Keras DataGenerators were used. DataGenerators allow for defining a python generator that provides batches of data to the fit method, instead of providing a fixed dataframe. This was useful when dealing with datasets consisting of time series of variable lengths and Fully Convolutional Networks. In such cases, the lengths of time series could be different for each batch. The preprocessing was designed differently given the model architecture. If the model expected time series of the same length `ConstantLengthDataGenerator` was used. For the Fully Convolutional Networks we implemented `VariableLengthDataGenerator`. The preprocessing steps for both generators are similar and will be described below:

### 3.3. NEURAL NETWORK ARCHITECTURES

- selecting time series to the batch. Number of series in a batch was equal to 32. The series were chosen so that the classes were uniformly distributed in the batch.
- normalizing the time series values to standard normal distribution. The normalization is done per one time series
- obtaining time series of the same length within the batch. For `ConstantLengthDataGenerator`. the length is constant and set a hyperparameter (256 by default). For `VariableLengthDataGenerator`. the length of the longest series in the batch is computed. The target length is equal to the length of the longest series or 3000 if it exceeds. In case of both generators if a time series is shorter than the target length, then the series is padded. If a time series is longer than the target length, then a random sub interval is selected.

### 3.3. Neural network architectures

The experiments were conducted on the Fully Convolutional Network. The network detailed architecture summary, including layer sizes is shown below.

Model: "FCN"

Layer (type)	Output Shape	Param #
input_1 (InputLayer)	[(32, None, 1)]	0
conv1d (Conv1D)	(32, None, 128)	1536
batch_normalization (Batch Normalization)	(32, None, 128)	512
activation (Activation)	(32, None, 128)	0
conv1d_1 (Conv1D)	(32, None, 256)	164096
batch_normalization_1 (Batch Normalization)	(32, None, 256)	1024

activation_1 (Activation)	(32, None, 256)	0
conv1d_2 (Conv1D)	(32, None, 128)	98432
batch_normalization_2 (Batch Normalization)	(32, None, 128)	512
activation_2 (Activation)	(32, None, 128)	0
global_average_pooling1d (GlobalAveragePooling1D)	(32, 128)	0
dense (Dense)	(32, 2)	258

=====

Total params: 266,370

Trainable params: 265,346

Non-trainable params: 1,024

The second parameter in output shape is always equal to **None**, because the length of the time series is not known in advance.

### 3.4. Multi-source transfer learning

As single source transfer learning for time series classification was covered by Fawaz et al. [6], this thesis considers mostly transfer learning with the usage of multiple source datasets. Two approaches of transfer learning were developed, named *baseline approach* and *ensemble approach*, both of them will be described below. Both of those approaches consider multiple (the number is arbitrary) source datasets, single target dataset, the same neural network architectures and hyperparameters.

#### 3.4.1. Baseline approach

The *baseline approach* is a simple extension to classical transfer learning approach used in deep learning. Usually in deep learning, the neural network is first trained on a large diverse

### 3.4. MULTI-SOURCE TRANSFER LEARNING

dataset (like ImageNet). The network is usually trained on the source dataset, then the first layers are used in the source classifier.

The idea here is similar, but the source dataset is created artificially from multiple datasets. As most of the datasets in the UCR archive are small, we concatenate them to obtain a bigger dataset. Suppose that we have  $N$  datasets  $(X_i^S, Y_i^S), i = 1, \dots, N$ . The target dataset  $(X^T, Y^T)$  is created as follows:  $X^T = \bigcup_{i=1}^N X_i^S$ ,

$Y = \{prefix_i(y) : y \in Y_i, i = 1, \dots, N\}$  where  $prefix_i(\cdot)$  is a function such that

$$\forall i \neq j \forall \hat{y} \in Y_i, \hat{y} \in Y_j \quad prefix_i(y) \neq prefix_j y$$

In our case, the prefix function given a class name (as a string), simply transforms the string by adding a prefix with the dataset name. This way the classes remain distinct after concatenation. The number of classes in the resulting dataset is equal to the sum of number of classes in each of the datasets. The examples that correspond to the same dataset should be intuitively similar to each other (analogue of the fine-grained classes in ImageNet) while examples from different datasets cover more latent space (analogue of the coarse-grained classes from ImageNet). The concatenation process is shown on figure 3.1

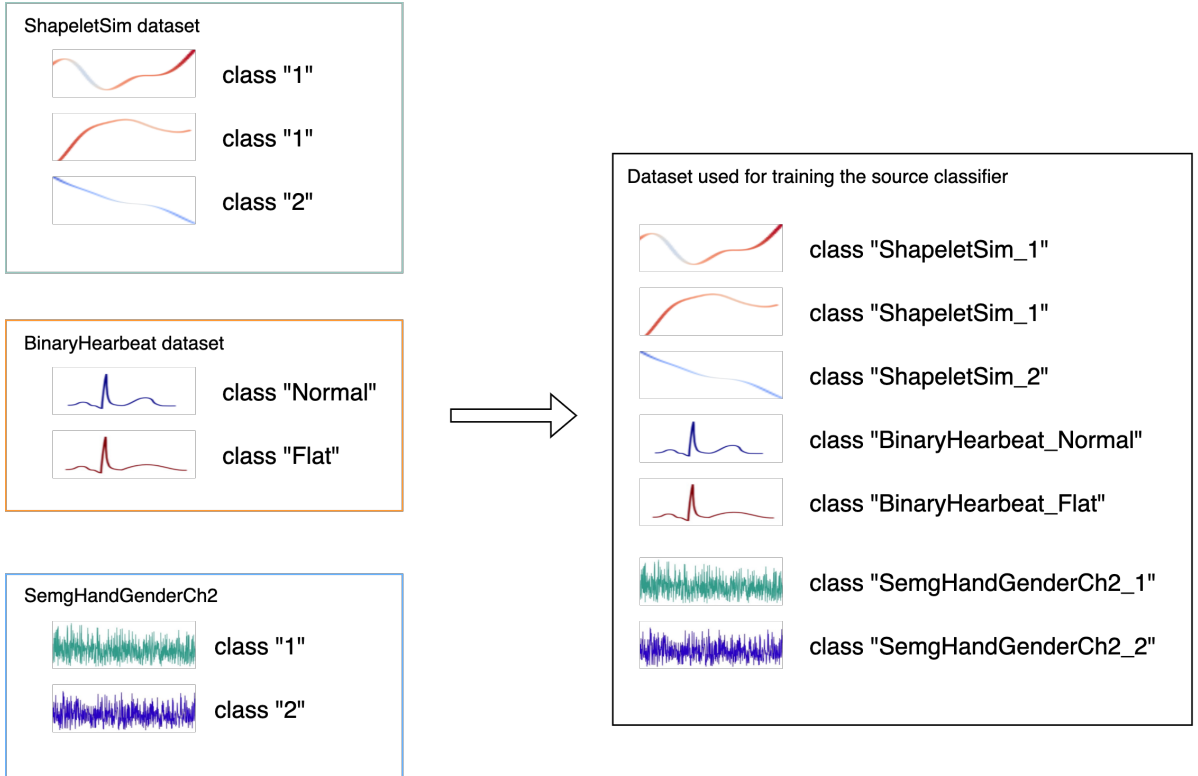


Figure 3.1: Illustration of preparing the data for the source classifier in the baseline approach.

The data batches that are feed to the model are processed as described in the previous chapter, including length unification within batch, normalization and random augmentation of 30% of the

observations.

After training the classifier on the artificial, concatenated dataset, we obtain a neural network model. The network (specifically the weights) is then used as a starting point to train on the target dataset. The target classifier is created by removing the last layer, which outputs the probabilities for each class. The layer is replaced by a new layer, with randomly initialized weights and output size equal to number of classes in the last layer. The whole network is then retrained on the target dataset. All weights are updated in the training procedure (freezing the weights is not used).

We compare the results by predicting the classes on the test split of the target dataset. The metrics used for comparison will be explained in detail in the Results chapter. We compare the metrics to results obtained by training a network of the same architecture but without transfer learning (with random weights).

### 3.5. Ensemble approach

The second approach proposed in this thesis is inspired by the success of ensemble models (e.g. HIVE-COTE) in time series classification. Small datasets sizes common in time series classification tend to increase the variance error of the classifiers. Therefore the ensemble method is often proposed to decrease the variance and increase accuracy of predictions. This intuition is resembled in this approach.

The previous study on transfer learning with a single source dataset [6] shown that though transfer learning generally increases the accuracy, for a lot of experiments the accuracy was decreased. If a dataset differed significantly from the target dataset, a classifier pretrained on such dataset wouldn't adapt to the target dataset, thus decreasing the accuracy. The same situation is probable with multi-source transfer learning, if part of the datasets are not suited for the particular target dataset. In order to mitigate the negative contribution, ensemble method is used.

Let's assume that we want to perform transfer learning given  $N$  source datasets  $(X_1^S, y_1^S), \dots, (X_N^S, y_N^S)$  and a target dataset  $(X^T, y^T)$ . The model is created from  $N$  source models. Model  $M_i$  is a classifier (e.g. a Fully Convolutional Network) trained on the dataset  $(X_i^S, y_i^S)$ . The hyperparameters used for experiments are described in section (TODO) and do not differ from the hyperparameters used for the *baseline approach*.

Let's recall the notion from the section 2.1.1 that a model  $M_i$  is a composition of layers,

$$M_i(X; \theta_1^i, \dots, \theta_M^i, \beta_1^i, \dots, \beta_M^i) = L_M^i(\dots L_2(L_1(X; \theta_1^i, \beta_1^i); \theta_2^i, \beta_2^i); \theta_M^i, \beta_M^i)$$



### 3.5. ENSEMBLE APPROACH

To shorten the notation, let's group the first layers together:

$$E(X; \Theta^i, B^i) = L_{M-1}^i(\dots L_2(L_1(X; \theta_1^i, \beta_1^i); \theta_2^i, \beta_2^i); \theta_{M-1}^i, \beta_{M-1}^i)$$

thus

$$M_i(X; \theta_1^i, \dots, \theta_M^i, \beta_1^i, \dots, \beta_M^i) = L_M^i(E(X; \Theta^i, B^i); \theta_M^i, \beta_M^i)$$

We assume that the architecture of subsequent layers except for the last one is the same for all models  $M_i$ ,  $i = 1, \dots, N$ . Therefore only the last layer has the superscript  $i$ . The last layer output size is equal to the number of classes in dataset  $i$ , therefore may be different for every model. We also assume that the first layers architecture is the same for all models. Models agnostic to the input size (like FCNs) are advantageous here, as their architecture does not depend on the input size. As for other models, we assume a fixed size of the input data achieved by preprocessing the input data. The weights  $\theta_k^i$  and biases  $\beta_k^i$  are obtained from training on dataset  $i$ . The target model  $M$  is created as follows:

$$M(X) = \frac{1}{N} \sum_{i=1}^N \hat{L}_M(E(X; \Theta^i, B^i); \hat{\theta}_M^i, \hat{\beta}_M^i)$$

The weights  $\theta_1^i, \dots, \theta_{M-1}^i$  and biases  $\beta_1^i, \dots, \beta_{M-1}^i$  are obtained from training the source models. The last layer  $\hat{L}_M$  is "new", meaning that the weights are randomly initialized. The new corresponding weights  $\hat{\theta}_M^i$  and biases  $\hat{\beta}_M^i$  have proper dimensionality (corresponding to number of classes in target dataset). This model is then retrained on the target dataset. All weights are updated in the training procedure (freezing the weights is not used).

A graphical example of the process can be found in figure. 3.2.

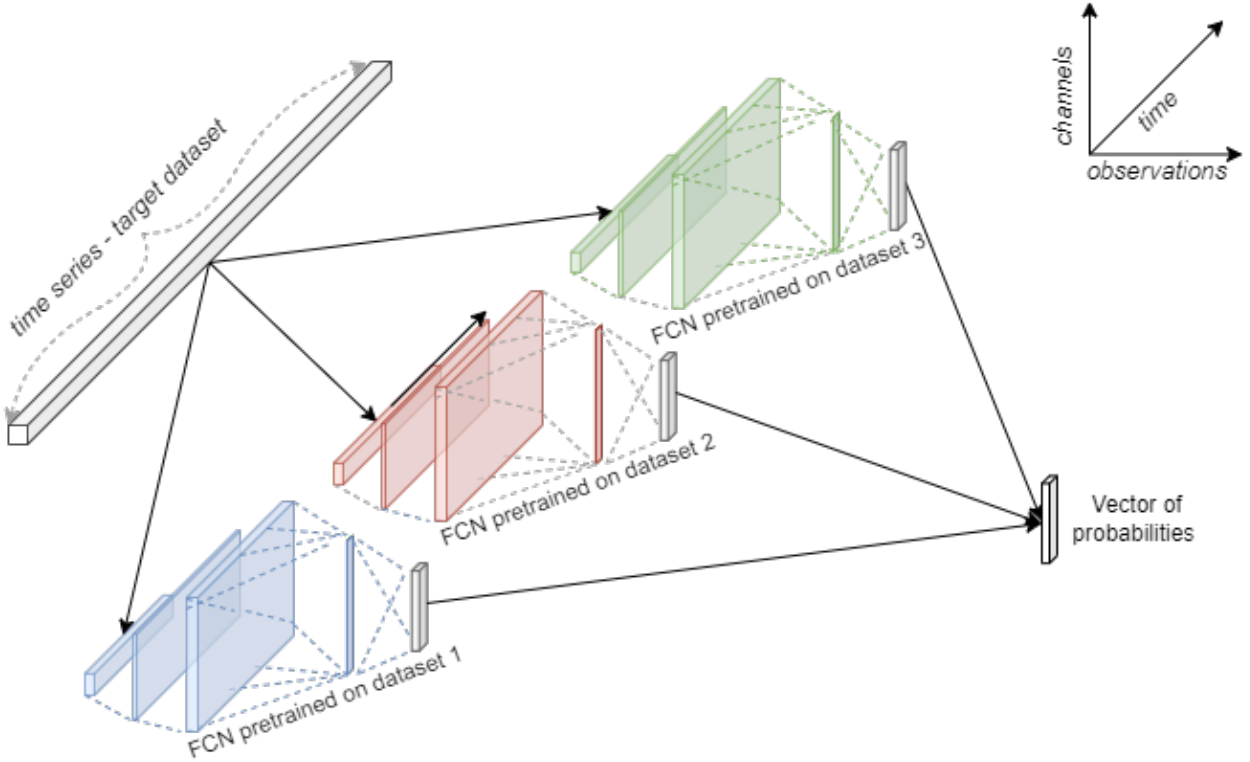


Figure 3.2: Illustration of the ensemble approach based on  $N = 3$  source models (Fully connected networks in this case). If a layer is coloured, it means that the weights in this layer are initialized by pre-training on one of the source datasets. The last layer in each model is new.

### 3.6. Selecting the source datasets for multi-source transfer learning

In both multi-source learning approaches described above (*baseline* and *ensemble* approaches) we don't impose any restrictions on the source datasets itself. Neither of the two approaches don't include any rules how the source dataset correspond to the target dataset. In the experiments conducted in this thesis, two approaches for picking the datasets are proposed. In the first approach, we pick randomly  $N$  datasets from a given category corresponding to the target dataset. To illustrate, let's say that the objective is to train a target classifier on the **InlineSkate** dataset, which belongs to the **Motion** category in the UCR archive. Then the source datasets are picked randomly from the **Motion** category, which could be **UWaveGestureLibraryX**, **AsphaltPavementType**, **ToeSegmentation2**, **CricketZ**, **UWaveGestureLibraryZ**.

The second approach described was based on a similarity measure. The inter-dataset similarity measure called DTW barycenter averaging was described in 2.4.4. The intuition behind the second approach is that usually transfer learning benefits if the domains are similar. To measure the similarity, the DBA distances were computed for all pairs of datasets in a given category.

### 3.6. SELECTING THE SOURCE DATASETS FOR MULTI-SOURCE TRANSFER LEARNING

Then, for a given target dataset,  $N$  datasets minimalizing the distance to the target dataset were used for pretraining the source network (either in *baseline* or *ensemble* approach). The dataset still belonged to the same UCR category as the target dataset.

## 4. Results

The results and metrics of classifiers obtained from training networks with transfer learning were compared to the results obtained without transfer learning. The performance of the target classifier in the ensemble approach was compared to an ensemble model with the same amount of component models in the ensemble.

### 4.1. Experiments

Each experiment recorded in MLFlow corresponds to one training setting. An example setting may be e.g.

- training a network using the *ensemble approach*, with 5 datasets picked using the DBA similarity measure
- training a network using the *baseline approach*, with 5 source datasets picked randomly
- training an ensemble of FCNs without transfer learning

Following the MLFlow convention, each experiments consisted of several runs, each correspond to one target dataset from the category **MOTION**. When comparing averaged results, then the result was averaged on all runs.

#### 4.1.1. Evaluation of the baseline approach

The *baseline* approach, in which a single FCN was pretrained and fine-tuned, was compared to results obtained from trained a single FCN model on the target dataset only. For this comparison, we used Fully Connected Network as the model and 5 source datasets selected to minimize the DBA distance to the target dataset. Figure 4.1 and shows ?? shows the accuracy and loss in each epoch, averaged over all experiments conducted.

#### 4.1. EXPERIMENTS

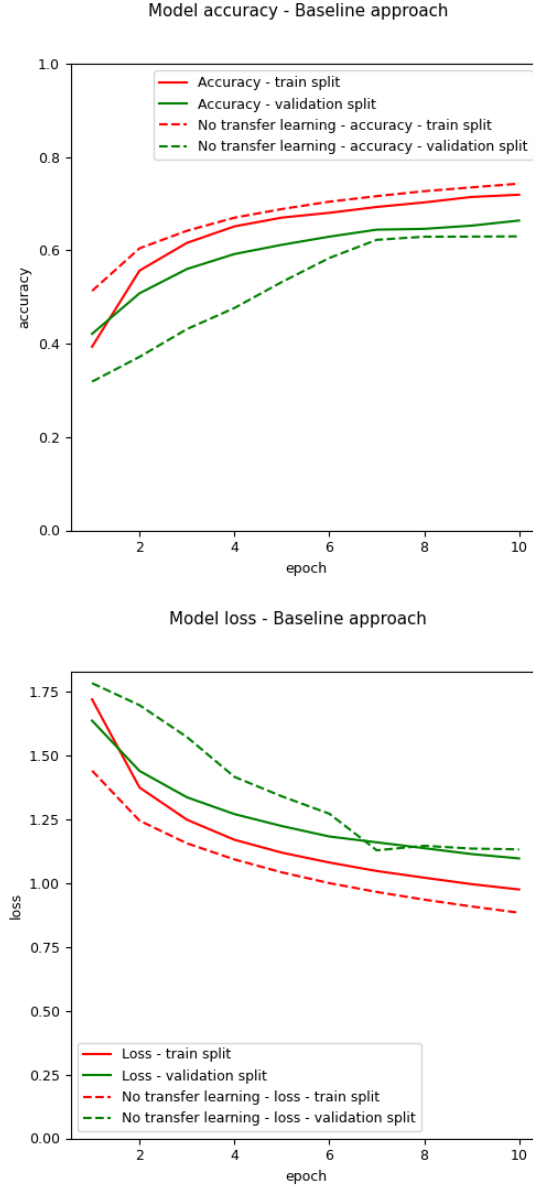


Figure 4.1: Averaged accuracy and loss obtained from transfer learning with the *baseline* approach, compared to a single FCN trained only on the target dataset.

The greatest improvement is visible in the first epochs of the training process, which suggest that the network uses the knowledge learned from the previous task. As the training continues, the advantage is more significant. On Figure 4.2 we show a comparison win/tie/loss diagram for fifth and tenth epoch. We also see that the gap between accuracy/loss on the validation and test split is smaller than without transfer learning. This can mean that the model is less prone to overfitting hence generalizes better, due to exposure to more diverse examples in the pre-training phase.

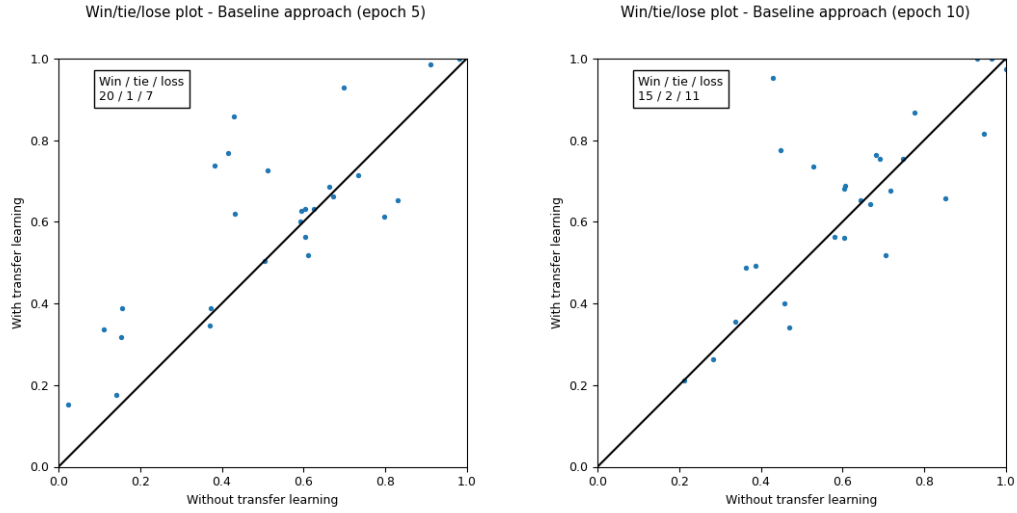


Figure 4.2: Win/tie/loss plot. This plot shows the validation accuracy in 5th and 10th epoch for each run (corresponding to a target dataset). The y axis corresponds to baseline approach, and x axis to approach without transfer learning. If a dot is above the diagonal, then it means that the transfer learning was beneficial.

## Bibliography

- [1] Anthony Bagnall, Aaron Bostrom, James Large, and Jason Lines. *The Great Time Series Classification Bake Off: An Experimental Evaluation of Recently Proposed Algorithms. Extended Version*. arXiv, 2016.
- [2] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate, 2014.
- [3] Zhicheng Cui, Wenlin Chen, and Yixin Chen. Multi-scale convolutional neural networks for time series classification, 2016.
- [4] Hoang Anh Dau, Anthony J. Bagnall, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, and Eamonn J. Keogh. The UCR time series archive. *CoRR*, abs/1810.07758, 2018.
- [5] Hoang Anh Dau, Eamonn Keogh, Kaveh Kamgar, Chin-Chia Michael Yeh, Yan Zhu, Shaghayegh Gharghabi, Chotirat Ann Ratanamahatana, Yanping, Bing Hu, Nurjahan Begum, Anthony Bagnall, Abdullah Mueen, Gustavo Batista, and Hexagon-ML. The ucr time series classification archive, October 2018. [https://www.cs.ucr.edu/~eamonn/time\\_series\\_data\\_2018/](https://www.cs.ucr.edu/~eamonn/time_series_data_2018/).
- [6] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Transfer learning for time series classification. In *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, dec 2018.
- [7] Hassan Ismail Fawaz, Germain Forestier, Jonathan Weber, Lhassane Idoumghar, and Pierre-Alain Muller. Deep learning for time series classification: a review. *Data Mining and Knowledge Discovery*, 33(4):917–963, mar 2019.
- [8] Minyoung Huh, Pulkit Agrawal, and Alexei A. Efros. *What makes ImageNet good for transfer learning?* arXiv, 2016.
- [9] Arthur Le Guennec, Simon Malinowski, and Romain Tavenard. Data Augmentation for Time Series Classification using Convolutional Neural Networks. In *ECML/PKDD Work-*

*shop on Advanced Analytics and Learning on Temporal Data*, Riva Del Garda, Italy, September 2016.

- [10] François Petitjean, Alain Ketterlin, and Pierre Gançarski. A global averaging method for dynamic time warping, with applications to clustering. *Pattern Recognition*, 44(3):678–693, 2011.
- [11] Joan Serrà, Santiago Pascual, and Alexandros Karatzoglou. Towards a universal neural network encoder for time series, 2018.
- [12] Chuanqi Tan, Fuchun Sun, Tao Kong, Wenchang Zhang, Chao Yang, and Chunfang Liu. A survey on deep transfer learning, 2018.
- [13] Zhiguang Wang, Weizhong Yan, and Tim Oates. Time series classification from scratch with deep neural networks: A strong baseline, 2016.
- [14] Karl Weiss, Taghi M. Khoshgoftaar, and DingDing Wang. *A survey of transfer learning*. Journal of Big Data, 2016.
- [15] Jinsung Yoon, Daniel Jarrett, and Mihaela van der Schaar. Time-series generative adversarial networks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [16] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. *A Comprehensive Survey on Transfer Learning*. arXiv, 2019.