



Intelligent Systems Project Report

Author:

Francesco Iemma

M.Sc. IN COMPUTER ENGINEERING

Academic Year 2020/21

Contents

1	Data Cleaning	3
1.1	Data Balancing	4
1.2	Features Selection	6
2	Neural Networks	7
2.1	Fitnet	7
2.2	RBF	7
3	Fuzzy Inference System	8
4	Convolutional Neural Networks	9

Introduction

The tasks performed in this project are the following:

- *3.1* "Design and develop two MLP artificial neural networks that accurately estimate a person's valence and arousal, respectively" and "two RBF networks that do the same thing as the MPLPs"
- *3.3* "Design and develop a fuzzy inference system to fix the deficiencies in the arousal dimension"
- *4.1* "Design and develop a convolutional neural network (CNN) that accurately classifies a person's emotion, based on facial expression."

The dataset at our disposal are two, one for tasks *3.1* and *3.3* and another one for task *4.1*. For what concern the first dataset, i.e. the one with biomedical signals for estimate arousal and valence, it is important to perform a cleaning of the data in order to obtain better performance for the neural networks that will be trained on it. This process, which is performed by the script `/matlab/data.m` is explained in the chapter 1.

After this chapter for each task is dedicated a chapter in which are explained the choices done and the results obtained in terms of performance.

Chapter 1

Data Cleaning

In this chapter we will see the data cleaning procedure performed in order to obtain better performance for the NNs. The steps done are:

- Remove non numeric values
- Remove outliers
- Balance the data among the different values of arousal and valence
- Features selection

All the procedure is contained in the file `/matlab/data.m`. The first two steps are performed thanks two matlab functions:

- `isinf(A)` that given a matrix returns a logic matrix is indicated if the correspondent element of the input matrix are infinite (1) or not (0)
- `rmoutliers(dataset, method)` that given a dataset remove the outliers found using the method specified in input that is 'median' by default (i.e. "Outliers are defined as elements more than three scaled MAD from the median. The scaled MAD is defined as $c \times \text{median}(\text{abs}(A - \text{median}(A)))$ ")

Then after the first two steps there are the most interesting part: data balancing and features selection.

1.1 Data Balancing

The dataset is composed by samples and each sample contains biomedical signals: to each set of biomedical signals (that we will call *features*) correspond a value for arousal and a value for valence. The possible values are 7 (1, $2.\bar{3}$, $3.\bar{6}$, 5, $6.\bar{3}$, $7.\bar{6}$, 9), thus we can divide the dataset according 7 class for arousal and valence. The distribution of the samples among the classes is represented in the histograms in figure 1.1 and 1.2.

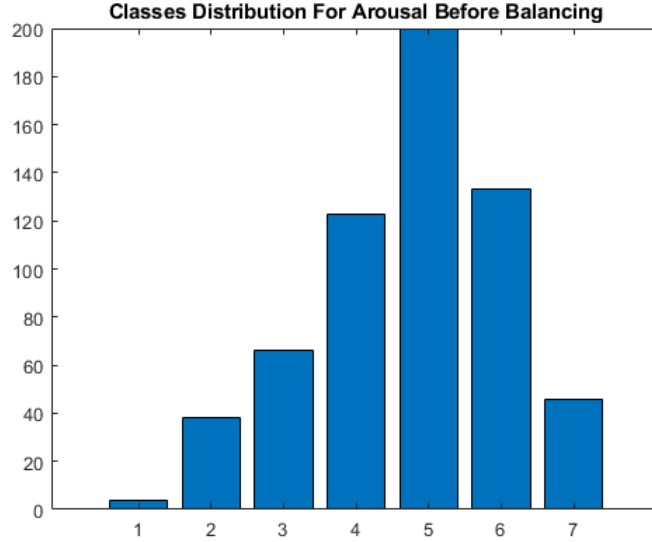


Figure 1.1: Classes Distribution For Arousal Before Balancing

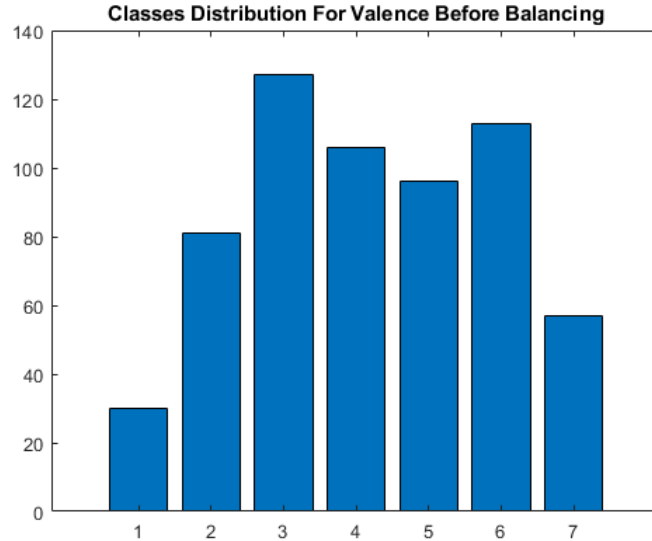


Figure 1.2: Classes Distribution For Valence Before Balancing

As we can see the samples are heavily unbalanced, for this reason an algorithm to balance the data has been used. The algorithm is based on the concepts of undersampling, oversampling and data augmentation.

The steps are the following:

1. I augment the samples that belong to the majority class of arousal and don't belong to the majority class of valence and viceversa (i.e. the samples that belong to the majority class of valence and don't belong to the minority class of arousal).

2. I remove the samples that belong to the majority class of arousal and don't belong to the minority class of valence and viceversa (i.e. the samples that belong to the majority class of valence and don't belong to the minority class of arousal).
3. I repeat steps 1 and 2 for a $n = 40$ (40 after some experiments this is the number that gives the best results) times and for each repetition I compute the new majority and minority class both for arousal and valence.
4. After the end of the repetitions I perform an undersampling on the first class because it is unbalanced for what concern the valence. Thus I remove some samples from this class, after some experiments removing 30 samples results in a balanced distribution for both arousal and valence.

At the end of this procedure the data are balanced as we can see in figure 1.3 and 1.4.

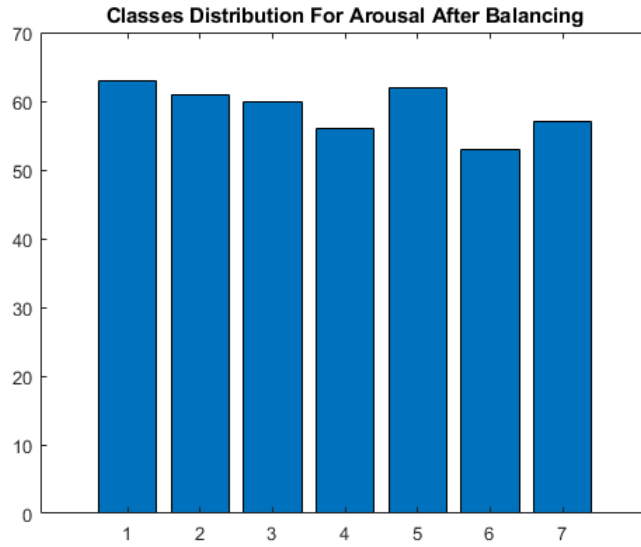


Figure 1.3: Classes Distribution For Arousal After Balancing

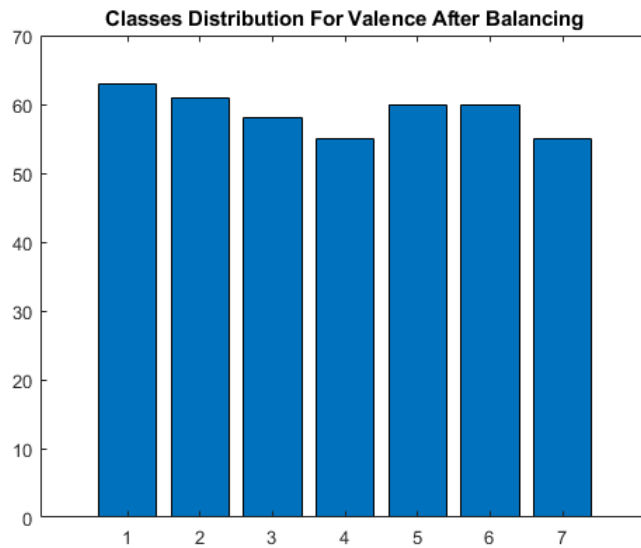


Figure 1.4: Classes Distribution For Valence After Balancing

1.2 Features Selection

Before the features selection is necessary to divide the data in two set: one for training and one for test. This is very important because if we use all the data to perform feature selection we have a bias because the test data have been already seen by the net.

Thus after the extraction of the holdout partition we perform x CAMBIARE times `sequentialfs` for arousal and x time for valence. Then we select the first `FEATURES_TO_SELECT` (constant set at the beginning of the script) features that appear most times in the different runs of `sequentialfs`, this operation is performed separately for arousal and valence.

At the end we save the data obtained into three `.mat` files:

- `/matlab/data/biomedical_signals/dataset_cleaned.mat`

It contains the entire dataset without infinite values, outliers and with balanced class distribution.

- `/matlab/data/biomedical_signals/training_data.mat`

It contains a `struct` with the training input (only the selected features) and the correspondent target output.

- `/matlab/data/biomedical_signals/test_data.mat`

It contains a `struct` with the test input (only the selected features) and the correspondent expected output.

Chapter 2

Neural Networks

In this chapter we will see two types of neural networks that resolve the same problem, that is to estimate the values of arousal and valence given a set of biomedical signals.

2.1 Fitnet

2.2 RBF

Chapter 3

Fuzzy Inference System

Chapter 4

Convolutional Neural Networks