# Using N-Grams of Quantized EEG Values for Happiness Detection

David Pinto[(✉)], Darnes Vilariño, Illiana Morales, Cristina Aguilar, and Mireya Tovar

Faculty of Computer Science Language and Knowledge Engineering Lab, Benemérita Universidad Autonóma de Puebla, Puebla, Mexico
{dpinto,darnes,mtovar}@cs.buap.mx, illiana.mrls.t@gmail.com,
crisaguilarc@hotmail.com
http://www.lke.buap.mx

**Abstract.** When applying classification methods for the automatic detection of happiness in human beings using electroencephalographic signals, the major research works in literature report the employment of power spectral density as the main feature. However, the aim of this paper is to explore wheter or not the use of N-grams of quantized EEG values as new features may help to improve the classification process. N-grams is a standard method of data representation in the area of natural language processing which usually reports good results. In this type of input data make sense to employ this kind of representation because the happiness signal is made up of a sequence of values which naturally matches the N-grams paradigm. The results obtained show that this kind of representation obtains better results than others reported in literature.

**Keywords:** EEG · N-grams · Happiness detection · Classification

## 1 Introduction

Among different emotional states the human being has, there is one defined by positive or pleasant emotion which ranges from content to intense joy named "happiness". The automatic detection of emotional states, in particular happiness, is the aim of this paper. Constructing a computational method able to recognize this emotional state in people surrounding us is an important part of human interaction, and also human-machine interaction. There have been different studies in literature for approaching the automatic detection of emotional states which can generally be categorized into three approaches, when the input signal is taken into account: (a) based on facial expression analysis, (b) based on electroencephalographic signals (EEG) and, (c) based on peripheral physiological signals [1]. The experiments carried out in this paper use EEG signals for

detecting whether or not a human being experiments an emotion of happiness, therefore, this paper can be categorized into the second approach.

There is, however, other types of categories we may use for those papers reporting emotional states detection in literature. Some papers differ from others because of the type of features employed in the signal representation. Some authors, for example, use EEG signals as input data considering Power Spectral Density (PSD) as the main feature with accuracies around 75 % and 65 % [2].

In [3], the authors propose a novel feature named "functional connectivity" which can be used together with other features in the task of automatic identification of emotional states. The authors report an accuracy between 0.4 and 0.65 for the experiments carried out. Even if there are other works in literature reporting results for the detection of emotional states by employing EEG signals, to the best of our knowledge works employing N-grams for the representation of data nearly have been reported in literature.

The protocol employed for acquisition is an issue that some research papers emphasize, for example, in [4], the authors have designed an acquisition protocol based on the recall of past emotional life episodes to acquire data from both peripheral and EEG signals. They report the performance of several classifiers for distinguishing between the three areas of the valence-arousal space, corresponding to negatively excited, positively excited, and calm-neutral states. The same authors propose an approach for affective representation of movie scenes based on the emotions that are actually felt by spectators [5].

The remaining of this paper is structured as follows. In Sect. 2 we present the concepts related with electroencephalography, in particular, we describe the wave patterns employed as input signals. The methodology proposed in this paper is given in Sect. 3. Here we describe the dataset, the quantization process, the data representation model and the algorithms for automatic identification of happiness. The obtained results are discussed in Sect. 4. Finally, in Sect. 5 the conclusions are given.

## 2   Electroencephalograpy

EEG refers to the recording of the brain's spontaneous electrical activity over a period of time [6]. Usually, this recording is done by using multiple electrodes which are placed on the scalp of a human being. The EEG is typically described in terms of rhythmic activity, which is divided into bands by frequency. It has been noted that these frequency bands have certain biological significance and distribution over the scalp.

In general, waveforms are subdivided into bandwidths known as *alpha*, *beta*, *theta*, and *delta* to describe the majority of the EEG signals used in clinical practice [7]. In the following subsection we briefly describe each one of these waveforms.

## 2.1   Wave Patterns

*Delta* is the frequency range up to 4 Hz. It tends to be the highest in amplitude and the slowest waves. This waveform is usually most prominent frontally in adults, whereas it is most prominent posteriorly in children. It is said that this waveform is normally seen in adults when they are in slow wave sleep. It is also seen in babies.

*Theta* is the frequency range from 4 Hz to 7 Hz. It is associated with some reactions as drowsiness or arousal in older children and adults, but it could be also seen when the person is in meditation [8]. When this waveform is presented in excess, it may represent abnormal activity. However, this range can be also be associated with reports of relaxed, meditative and creative states.

*Alpha* is the frequency range from 7 Hz to 14 Hz. Hans Berger named the first rhythmic EEG activity he saw as the "alpha wave" [9]. This was the "posterior basic rhythm" seen in the posterior regions of the head on both sides, higher in amplitude on the dominant side. It emerges with closing of the eyes and with relaxation, and attenuates with eye opening or mental exertion. The alpha frequency range in young children is slower than 8 Hz and, therefore, it is technically in the theta range.

*Beta* is the frequency range from 15 Hz to about 30 Hz. It is associated to motor behavior and it is normally attenuated during active movements [10]. When this signal presents low amplitude with multiple and varying frequencies, it is frequently associated with active, busy or anxious thinking and active concentration. Rhythmic beta with a dominant set of frequencies is associated with various pathologies and drug effects. This is the dominant waveform in persons who are anxious or alert.

In Table 1 we show the frequency ranges for each one of the waveforms aforementioned.

**Table 1.** EEG Frequency bands

| Band | Frequency (Hz) |
|------|----------------|
| Delta | < 4 |
| Theta | = 4 and < 8 |
| Alpha | = 8 and < 14 |
| Beta | = 14 |

## 3   Methodology Proposed

In this section we describe the design cycle we have proposed for the analysis of EEG signals with the purpose of identifying whether a given sequence of waveforms expresses an emotion of happiness.

### 3.1   Dataset Construction

In order to collect a dataset for the experiments, we have created an ad-hoc software which is able to record the brain signal of a given person while he or she is observing a video containing scenes which trigger an emotion of happiness. The video has previously been annotated in order to determine which parts are associated to happiness and the parts of the video that are not. In this paper we are not interested in annotating other type of emotions, because we are only determining those sequences associated to happiness, whereas the remaining sequences are just annotated as unhappiness.

We employed three different videos with a length time of 4.2, 4.1 and 6.3 min, respectively. The three videos were shown to 20 different persons in order to construct the dataset. The brain signals of each person were parallelly recorded while they were observing the video. The adquisition software developed was able to automatically annotate the output signal with the corresponding tag associated to happiness emotion (H) or another tag associated to unhappiness emotion (N) for each one of the four waveforms: Delta, Theta, Alpha and Beta.

### 3.2   Quantization

Quantization is the process of mapping a large set of input values to a (countable) smaller set. This process may help to improve the performance of different tasks because all those similar values are grouped together in a single value which represent to all of them. The quantization process we have applied to the EEG signals are based on the mean and standard deviation values obtained from the EEG signals. In particular, we have obtained these two values ($\mu$ and $\sigma$) for each one of the four channels Delta, Theta, Alpha and Beta, and, therefore, the original signal may be grouped together into three different clusters, LOW, MEDIUM and HIGH, which may be better as discrete values used instead of the continuous ones. The limits of the clusters were found considered that the distribution of the waveforms for the person may be modeled as a normal distribution. In Table 2 we may see the ranges used for the quantization process.

**Table 2.** Quantization ranges

| Discrete value | Symbol | Range |
|---|---|---|
| LOW | L | EEG_Signal $< \mu - \sigma$ |
| MEDIUM | M | $\mu - \sigma \leq$ EEG_Signal $< \mu + \sigma$ |
| HIGH | H | EEG_Signal $\geq \mu + \sigma$ |

Although, these quantization thresholds have performed well in the experiments carried out, we consider that more investigation needs to be done with respect to the task of determining the optimal thresholds, so as to analyze whether or not, another probabilistic distribution should be employed instead of the normal one.

In order to distinguish the four waveforms in the discretization process, we have prefixed each quantized EEG signal (represented with the symbols: L, M and H) with the first letter of the waveforms, i.e., we used the letter "A" for Alpha, "D" for Delta and so on. In Table 3 we may observe a sample of the quantized sequences obtained. Each value of the sequence represents one second of measure of the person human brain. Each sequence has been annotated with a tag indicating whether or not the sequence is associated with an emotion of happiness.

**Table 3.** Sample of quantized sequence of EEG signals. The $H$ tag means "Happiness", whereas the $U$ tag means "Unhappiness"

| Sequence of quantized EEG signals | Tag |
|---|---|
| AL AL AL AL AL AL AL AL | U |
| DL DL DL DL DM DL DL DL | U |
| TL TL TL TL TL TL TL TL | U |
| BL BL BL BL BL BL BL BL | U |
| AL AL AM | H |
| DL DM DM | H |
| TL TM TM | H |
| BL BL BM | H |
| AL AL AL AL AL AL AL AL | U |
| DL DL DL DL DL DL DL DL | U |
| TL TL TL TL TL TL TL TL | U |
| BL BL BL BL BL BL BL BL | U |
| AM AH AL AM AL AL AL | H |
| DL DM DL DL DL DL DL | H |
| TM TM TL TL TL TL TL | H |
| BM BL BL BL BL BL BL | H |
| AL AL AL AL AH AM | U |
| DL DL DL DL DM DM | U |
| TL TL TL TL TM TM | U |
| BL BL BL BL BM BM | U |

## 3.3   Data Representation

The quantization process previously applied allow us to have sequences of quantized values associated to a given tag. In the particular area of natural language processing, we describe this sequence of values as a string. Each string needs to be correctly represented in order to apply machine learning methods. The

feature extraction process is then applied by using the N-grams representation technique. Each string is split out into sequences of $n = 2, 3$ values, calculating the frequency of each N-gram. Since, the frequency is not sufficient for determining the degree of discrimination each N-gram has, we have employed a technique of term weighting known as TF-IDF in which each N-gram is weighted in terms of its frequency in the string and also proportional to the inverse document frequency (the number of strings containing the N-gram). The complete dataset is made up of documents manually annotated and represented by N-grams weighted with the TF-IDF schema.

### 3.4   Classification Model

The supervised machine learning techniques are able to learn the human process of identifying emotions based on features fed in the classifier by means of the manually annotated corpus.

We have selected one learning algorithm from four different types of classifiers: Bayes, Lazy, Functions and Trees in order to investigate the one that performs better in the particular task of automatic identification of happiness. The following four learning algorithms were chosen:

**NaïveBayes:** This is the standard probabilistic Naïve Bayes classifier.
**K-Star:** This is the $k$-nearest neighbor classifier with a generalized distance function.
**SMO:** This is a sequential minimal optimization algorithm for support vector classification.
**J48:** This is the C4.5 decision tree learner which implements the revision 8 of C4.5.

In the following subsection we describe the measures we employed for evaluating the performance of the experiments carried out.

### 3.5   Evaluation Measures

In order to evaluate the quality of the results obtained, we have used the following standard measures for the evaluation: Precision, Recall and $F$-Measure [11].

The measures employed make use of a set of values calculated when the classification process is carried out. The terms "true positive - (TP)", "true negative - (TN)", "false positive - (FP)", and "false negative - (FN)" compare the results obtained by the classifier under test set with a given gold standard wich is usually obtained by external judgments (manual annotated data). The terms "positive" and "negative" refer to the classifier's prediction, and the terms "true" and "false" refer to whether that prediction corresponds to the external judgment.

Thus, the Precision and Recall is calculated as follows:

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

The $F$-Measure combines Precision and Recall as the harmonic mean of these two values. The traditional F-measure or balanced F-score is calculated as follows:

$$F - Measure = \frac{Precision * Recall}{Precision + Recall}$$

In this paper we also use the Accuracy measure for reporting the results obtained. This value is calculated as follows:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

## 4   Experimental Results

In this section we are presenting the accuracy obtained by each classifier when attempting to identify whether or not a sequence of EEG signals correspond to a signal of human happiness.

Tables 4, 5, 6 and 7, show the detailed accuracy by class using the KStar, Naïve Bayes, SMO and J48 supervised classifiers, respectively. As can be seen, in all the cases the identification of Class 1 (when the set of N-grams of quantized EEG signals is associated to a human emotion of happiness) obtained a better performance than the identification of Class 2 (when the set of features does not correspond to a human emotion of happiness). Even if the difference is not so significant, this issue is important. As future work, we need to provide better features for improving the results obtained. On the one hand, the KStar classifier was the one that obtained the worst results with a weighted average $F$-Measure of 0.745. On the other hand, the other three classifiers obtained similar results with a weighted average $F$-Measure of around 0.79.

**Table 4.** Detailed accuracy by class using the KStar classifier

| Class | Precision | Recall | F-Measure |
|---|---|---|---|
| Class 1 (Happiness) | 0.667 | 0.786 | 0.721 |
| Class 2 (¬Happiness) | 0.760 | 0.633 | 0.691 |
| Weighted Avg | 0.715 | 0.707 | 0.745 |

**Table 5.** Detailed accuracy by class using the Naïve Bayes classifier

| Class | Precision | Recall | F-Measure |
|---|---|---|---|
| Class 1 (Happiness) | 0.742 | 0.875 | 0.803 |
| Class 2 (¬Happiness) | 0.860 | 0.717 | 0.782 |
| Weighted Avg | 0.803 | 0.793 | 0.792 |

**Table 6.** Detailed accuracy by class using the SMO classifier

| Class | Precision | Recall | F-Measure |
|---|---|---|---|
| Class 1 (Happiness) | 0.735 | 0.893 | 0.806 |
| Class 2 (¬Happiness) | 0.875 | 0.700 | 0.778 |
| Weighted Avg | 0.808 | 0.793 | 0.792 |

**Table 7.** Detailed accuracy by class using the J48 classifier

| Class | Precision | Recall | F-Measure |
|---|---|---|---|
| Class 1 (Happiness) | 0.732 | 0.929 | 0.819 |
| Class 2 (¬Happiness) | 0.911 | 0.683 | 0.781 |
| Weighted Avg | 0.825 | 0.802 | 0.799 |

**Table 8.** Percentage of correctly vs. incorrectly instances classified

| Classifier | Type | Correct (%) | Incorrect (%) |
|---|---|---|---|
| K-Star | Lazy | 70.69 | 29.31 |
| Naïve Bayes | Bayes | 79.31 | 20.68 |
| SMO | Functions | 79.31 | 20.68 |
| J48 | Trees | 80.17 | 19.82 |

In Table 8 we show the percentage of instances classified correctly and incorrectly. Basically, this table summarize the weighted average results of the previously shown result tables. Actually, we have included a graph (see Fig. 1) with the aim of showing the summarized average results.

As can be seen, the J48 classifier performed similar to SMO and Naïve Bayes, with an accuracy of 80.17 %. Also, it is noticeable that these three classifiers outperformed the K-Star classifier.

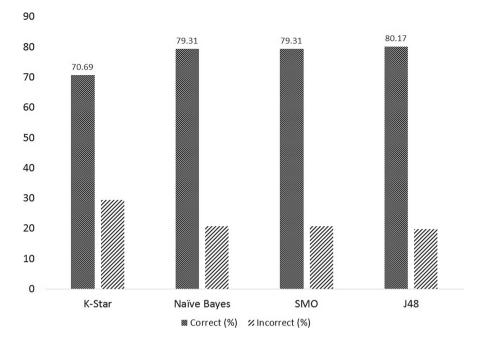**Fig. 1.** Comparison among the different classifiers employed in the experiments

## 5    Conclusions

In this paper we have presented a novel representation based on N-grams of quantized EEG values which obtained good results for the particular task of automatic happiness detection. The obtained accuracy results are up to 80 %, and the F-measure is about 0.8. Unfortunately, the obtained results can not be directly compared with those previously reported in the literature because we are not using the same datasets, nor the same type of sensors. A fair comparison among different methods should take into account that the experiments should be executed in very similar conditions.

The employment of N-grams is capturing information about a sequence of signals with acceptable performance, however, we consider we could improve the results obtained up to now, by proposing a different quantization method or even better, by increasing the number of samples used in the training phase.

In this paper we have considered that the data are distributed according to the normal distribution, a very common continuous probability distribution. However, we need still to investigate whether or not this assumption is correct for the quantization process. Instead, we could employ another probability distribution such as the gamma or the poison one.

# References

1. Wang, X.-W., Dan Nie, B.L.L.: Emotional state classification from EEG data using machine learning approach. Neurocomputing **129**(1), 94–106 (2014)
2. Jatupaiboon, N., Pan-ngum, S., Israsena, P.: Real-time EEG-based happiness detection system. Sci. World J. **2013**, 52–61 (2013)
3. Lee, Y.-Y., Hsieh, S.: Classifying different emotional states by means of EEG-basedfunctional connectivity patterns. PLoS ONE **9**(4), e95415 (2014)
4. Chanel, G., Kierkels, J.J., Soleymani, M., Pun, T.: Short-term emotion assessment in a recall paradigm. Int. J. Hum. Comput. Stud. **67**(8), 607–627 (2009)
5. Soleymani, M., Chanel, G., Kierkels, J.J.M., Pun, T.: Affective characterization of movie scenes based on content analysis and physiological changes. Int. J. Semant. Comput. **3**(2), 235–254 (2009)
6. Niedermeyer, E., da Silva, F.L.: Electroencephalography: Basic Principles, Clinical Applications, and Related Fields, 5th edn. Lippincott Williams & Wilkins, Baltimore (2004)
7. Tatum, W.O.: Ellen r. grass lecture: Extraordinary EEG. Neurodiagnostic J. **54**, 3–21 (2014)
8. Cahn, B.R., Polich, J.: Meditation states and traits: EEG, ERP, and neuroimaging studies. Psychol. Bull. **132**(2), 180–211 (2006)
9. Millet, D.: The origins of EEG. In: Seventh Annual Meeting of the International Society for the History of the Neurosciences (ISHN), Los Angeles, California, USA Department of Neurology, UCLA Medical Center (2004)
10. Pfurtscheller, G., de Lopes, F.H.: Event-related EEG/MEG synchronization and desynchronization:basic principles. Clin. Neurophysiol. **110**(11), 1842–1857 (1999)
11. Fawcett, T.: ROC graphs: Notes and practical considerations for researchers. Technical report, HP Labs (2004)