

# DYNAMIC SEQUENCE-BASED LEARNING APPROACHES ON EMOTION RECOGNITION SYSTEMS

Imen Trabelsi

UR Sciences and Technologies  
of Image and Telecommunications  
University of Sfax, Tunisia  
Email:imen.trabelsi@enit.rnu.tn

Mohammed Salim Bouhlel

UR Sciences and Technologies  
of Image and Telecommunications  
University of Sfax, Tunisia  
Email:medsalim.bouhlel@enis.rnu.tn

**Abstract**—Emotion recognition, especially from speech, is one of the key steps towards building an effective human-machine interaction. The recent progress from speaker recognition filed opens up a new area of exploration in emotion recognition. In fact, speaker recognition and emotion recognition have been always considered as separate systems operating independently of each other. However, due their their similarities, adapting any general emotional recognition system to the advanced Kernel-based machine learning algorithms from the field of speaker recognition, we argue that it is more efficient in terms of recognition performance. The focus in this paper is specially laid on building a range of dynamic sequence-based learning approaches including the KL divergence kernel, GUMI kernel and GLDS kernel. Extensive computer simulations were conducted on the well-established real-life speech dataset (IEMOCAP) and the acted Berlin emotional speech dataset (Emo-Db).

**Index Terms**—sequence kernel, speech emotion recognition, GMM, SVM, MFCC

## I. INTRODUCTION

Research on recognizing emotions is an important research topic in the area of human machine-interaction. Recently, speech emotion recognition, which aims to analyze the emotion states through speech signals, has been attracting increasing attention. Speech emotion recognition has several applications in day-to-day life. For example, agent-customer interactions can use emotion recognition systems to evaluate specific aspects of customer service quality, as well as customers preferences, likes and dislikes herm. Emotion recognition systems may be used in an on-board car driving system [14], where the drivers behavior is influenced by the situation (e.g., traffic jam) and information about where the emotional state, in turn, has an impact on the driving behavior and may be used to keep him alert during driving. Other tasks that rely on human behavior analysis, such as in therapeutic settings or on intimate relationships [12] that can also benefit from robust emotion recognition. Increasingly, intelligent E-tutoring applications are gradually enhanced with emotional awarness capabilities [11]. Such systems have to recognize and analyze emotional information of students' learning performance to foster interactions and positive evaluations. Analysis of Emotion in Call

centers conversations is also a very important problem from a business perspective and helps to improve the quality of service of a call attendant[13]. An important issue to be considered in the evaluation of an emotional speech systems is the choice of emotional corpus. The existing emotional databases could be divided into three classes: simulated, elicited and spontaneous corpus [10]. Simulated speech is acted speech, which is expressed by a professional actor/actress in a deliberated manner. A common method is to employ actors in a laboratory environment, where the utterances are linguistically and phonetically predefined as done in the Berlin Database of Emotional Speech. Therefore the lacking naturalness of the emotional is often criticized in literature. The earliest research on emotion recognition starts with acted emotional data, and then extends to elicited data, which are more natural than their simulated counterparts. Elicited (or induced) speech corpora are collected by simulating the artificial emotional situation, in an interaction between a speaker and a machine or in a Wizardof-Oz (WOZ) technique, without the knowledge of the speaker. The elicited emotion is the natural react to different contextual situations. However, the induced emotions are often mild, as if there were an inverse relation between the strength of the induction and the unethical value. Recently, the demanding for real application forces the research shift to authentic and spontaneous emotions. Natural emotional speech corpus is the recording in a real-life situations, without any control, e.g., to record call center dialogues or talk-shows, recordings during abnormal conditions, a dialogue between patient and a doctor and so on. The challenge in the application of in vivo methods such as recordings everyday speech is to find a wide emotion base in this category and also to achieve a corpus of good technical quality especially without background noise outside a laboratory. There are also some legal issues, such as privacy and legal copyright problem. Fig. 1 shows a generic speech emotion recognition system that will be used in this paper. The first step is to extract feature vectors from the speech signals (reference and test). Once this is done, there are many approaches to build the emotion model. The majority of speech emotion recognition systems in use

today are based on statistical classification methods. Many of the recent studies use generative machines learning such as Gaussian Mixture Model (GMM). Nevertheless, discriminant machines learning such as Support Vector Machines (SVM) can improve the accuracies better than generative classifiers. SVM has recently shown to be powerful in many speech based classification problems [1], [4], [6]. We show that the same performances can be achieved in the field of emotion classification via sequence kernels. A sequence kernel compares emotional utterances as entire sequence rather than a probability calculated at the frame level

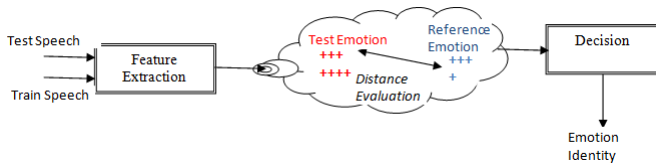


Fig. 1. Components of a typical speech emotion recognition system.

Kernel functions can be divided into two groups. The first one consists in training SVMs directly on the acoustics vectors which characterize the emotional utterances (frame level approach). The size of the set of feature vectors depends on the duration of the utterance and the number of the extracted features. The SVM with standard kernels like linear kernel, polynomial kernel and gaussian kernel cannot handle such varying length patterns. In addition, such a frame-level discriminative approach gives poor performance. Another approach consists in applying SVM as an inner product between the mapping of two sequences in an infinite feature space. In fact, the goal is to minimize classification errors on sequences, not on speech frames. That is why a sequence-based learning approach seems more appropriate. The kernels designed for varying length patterns are referred to as dynamic kernels. A number of dynamic sequence kernels have been developed recently, primarily targeted to speaker classification tasks. In this paper, we evaluate the generalized linear discriminant sequence kernel[15], the Kullback-Leibler divergence kernel[17], and the GMM-UBM mean kernel [16]. The remainder of the paper is organized as follows. In section 2, we provide details of the SVM employed. We introduce the various sequence kernels (outlined in Section 3 and 4). We then analyze the results comparing their individual performances. Finally, Section 6 gives a summary.

## II. SUPPORT VECTOR MACHINES

Support vector machine (SVM) is a powerful discriminative classifier that has been recently adopted in speaker/emotion recognition, and it has also been successfully combined with GMM to increase accuracy [1]. The SVM proposed by Vapnik [3] has been studied extensively for classification, regression and density estimation. This linear two-class classifier works by embedding the data into a Hilbert space(feature space), and searching a linear separator (hyperplane) in this space (shown in two feature dimensions in fig 2). The classifier equation is

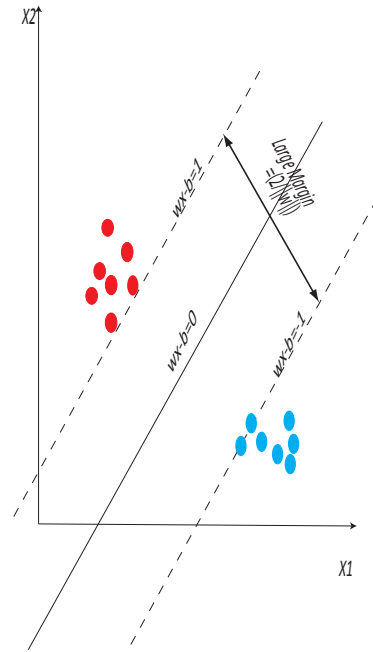


Fig. 2. Maximum-margin hyperplane.

given as follows:

$$f(x) = \sum_{i=1}^N \alpha_i y_i k(x, x_i) + d \quad (1)$$

where  $x_i$  are the support vectors chosen from training data via an optimization process, the target class label  $y_i \in -1, +1$

and  $N$  denotes the number of support vectors and  $d$  is a (learned) constant. The kernel function is constrained to satisfy the Mercers conditions, so that they can be expressed as  $K(x, x_i)$ , the kernel function that fulfills the Mercer conditions [?]:

$$K(x, y) = \phi(x)^t \phi(y) \quad (2)$$

where  $\phi(x), \phi(y)$  are a mapping from the input space to a high-dimensional space. SVMs have been used effectively in emotion classification [1]. A detailed description about SVM classifiers and how to construct a multi-class SVM classifier can be found in [3], [7]. The SVMs used in this study were implemented using LibSVM Toolkit [8].

## III. GMM SUPERVECTOR AND DYNAMIC SEQUENCE KERNEL FUNCTIONS

One of the issues in emotion recognition is how to represent utterances that, in general, have a varying number of feature vectors. The main idea of this proposed approach is to use a new feature representation based on GMM to construct the input vectors to train the SVM. This leads to sequence kernel SVM, where the utterances with variable number of feature vectors are mapped to a fixed-length vector, a so-called supervector. This supervector refers to combining many smaller-dimensional vectors into a higher and a fixed dimensional vector. It is important that the supervectors of different

utterances arise from a common coordinate system such as being adapted from a universal background model, or being generated using a fixed polynomial basis. For supervectors adapted from a universal background model (UBM), the first step is to compute this UBM representing a general structure of the underlying feature space of speech signals. This UBM is modeled as a big GMM model. Parameters for the UBM are trained using the EM algorithm which iteratively refines the GMM parameters to monotonically increase the likelihood of the estimated model for the observed feature vectors. After computing the UBM, the sequence of feature vectors extracted from the utterance are adapted (means only) from the UBM using A Posteriori (MAP) estimation [?] to produce emotion specific GMM given by the following equation.

$$p(x) = \sum_{i=1}^M w_i P(x; \mu_i, \Sigma_i) \quad (3)$$

where  $\mu_i, \Sigma_i$  denotes the mean and covariance of the  $i$ th multivariate Gaussian component  $P(x)$  and  $M$  is the total number of GMM components. From this GMM, the mean vectors are extracted and concatenated -typically normalized by the corresponding standard deviation of each of the Gaussian mixtures in the adapted GMM model to form the set GMM supervectors. The construction process of supervectors can be summarized as follows:

- A GMM model with  $M$  Gaussians is obtained via MAP adaptation from the UBM
- A kernel function is used to transform parameters of each GMM component to a fixed length vector. The vector corresponding to the  $i^{th}$  GMM component constitutes the  $i^{th}$  subvector of a supervector.
- All the subvectors are concatenated to obtain a high-dimensional vector.

In other words, this process transforms variable length utterances to a unique fixed-size vector which carries emotional information. This representation is in conformity with Eq. (2) where two arbitrary utterances  $x$  and  $y$  from the input space can be compared in the supervector space using the relation  $K(x, y) = \phi(x)^t \phi(y)$ , where  $K$  is the kernel function and  $\phi(x), \phi(y)$  are the supervectors obtained from utterances  $x$  and  $y$ , respectively.

In our experience, we have used four different kernel functions, The two first ones correspond to the Kullback-leibler kernel (the Linear GMM Supervector Kernel and the Non Linear GMM Supervector Kernel). The third one is the GUMI Kernel, The last one is the GLDS kernel.

#### IV. KULLBACK-LEIBLER KERNEL

##### A. Linear GMM Supervector Kernel

The kernel is based on a linear approximation of the Kullback Leibler (KL) distance between two gaussian distributions. Hence, KL kernel measures the distance between two supervectors. This distance corresponds to the Euclidean distance between scaled supervectors  $\lambda x$ , and  $\lambda y$ .

$$D^2(\lambda x, \lambda y) = \sum_{i=1}^N w_i (\mu_i^a - \mu_i^b) \Sigma_i^{-1} (\mu_i^a - \mu_i^b)^t \quad (4)$$

Where  $w_i$  and  $\Sigma_i^{-1}$  are the  $i$ th UBM mixture weights and diagonal covariance matrix,  $\mu_i$  correspond to the mean of the Gaussian  $i$  of the GMM emotional model.

The linear kernel is defined as the corresponding inner product:

$$K^{lin}(\lambda x, \lambda y) = \sum_{i=1}^N (\sqrt{w_i \Sigma_i^{-1/2}} \mu_i^a) (\sqrt{w_i \Sigma_i^{-1/2}} \mu_i^b)^t \quad (5)$$

##### B. Non Linear GMM Supervector Kernel

The non linear kernel is a Gaussian kernel defined on the GMM supervector space. The kernel function is expressed as an exponential function of distance  $D$ . The non linear GMM supervector kernel is obtained as:

$$K^{NonLin}(\lambda x, \lambda y) = e^{-D^2(\lambda x, \lambda y)} \quad (6)$$

#### V. THE GMM-UBM MEAN INTERVAL KERNEL

An alternative measure of separability between two probability distributions is the Bhattacharyya affinity measure. The GUMI (GMM-UBM mean interval) kernel is derived from the Bhattacharyya distance between  $p(x)$  and  $g(x)$  defined over  $x$  is given by:

$$\mathfrak{B}(p(x) \parallel g(x)) = \int_{-\infty}^{+\infty} \sqrt{p(x)} \sqrt{g(x)} dx \quad (7)$$

where  $p(x)$  and  $g(x)$  are two gaussian distributions. The GUMI kernel is obtained from the Bhattacharyya mean distance between a GMM and an UBM. This kernel is defined as the corresponding inner product:

$$K^{Gumi}(\lambda x, \lambda y) = \phi^{Gumi}(\lambda x)^t \phi^{Gumi}(\lambda y) = \sum_{i=1}^N (\mu_i^x - \mu_i)^t \sigma_i^{-1} (\mu_i^y - \mu_i) \quad (8)$$

#### VI. GENERALIZED LINEAR DISCRIMINANT SEQUENCE KERNEL

One of the simplest SVM kernel is the Generalized linear discriminant sequence (GLDS) kernel. This kernel function creates the supervectors not from a universal background model but by explicit mapping into kernel feature space using a fixed polynomial basis. The maximum degree of polynomials  $p$ , is a user-supplied parameter. Given a sequence of cepstral features  $x = (x_1, x_2, \dots, x_m)$ , the mapping function  $\phi^{GLDS}$  is expressed as:

$$\phi^{GLDS}(\lambda x) = \frac{1}{M} \sum_{m=1}^M \varphi(\lambda x_m). \quad (9)$$

The GLDS kernel between two examples  $X = x_1, x_2, \dots, x_m$  and  $Y = y_1, y_2, \dots, y_n$  is given as:

$$K^{GLDS}(\lambda x, \lambda y) = \left( \frac{1}{M} \sum_{m=1}^M \varphi(\lambda x_m) \right) S^{-1} \left( \frac{1}{N} \sum_{n=1}^N \varphi(\lambda y_n) \right). \quad (10)$$

where  $S = \frac{1}{T}R^T R$ ,  $T$  the total number of feature vectors from all the examples in the training data set and  $R$  is the matrix whose rows are the polynomial expansions of the feature vectors in the training set. In our implementation, the correlation matrix  $S$  is chosen to be diagonal. We estimate it using all the polynomial supervectors of the training set, as given in:

$$K^{GLDS}(\lambda x, \lambda y) = \frac{1}{M \sum_{m=1}^M S^{-\frac{1}{2}} \varphi(\lambda x_m) + (\frac{1}{N} \sum_{n=1}^N S^{-\frac{1}{2}} \varphi(\lambda y_n))}. \quad (11)$$

## VII. EMOTIONAL DATABASES

To generalize our study, two publicly available emotional databases, the Emo-Db database and the IEMOCAP database, are used in our experiments to validate the performance of the proposed emotion-recognition systems.

### A. EMO-Db database

EMO-Db is recorded by speech work-group led in the anechoic chamber of the Technical University in Berlin. It is a simulated open source speech database. In this database, ten professional native German actors (5 female and 5 male) simulated 7 emotions, producing 10 utterances. 5 utterances are short, while the remaining 5 are long. The emotions are: anger, boredom, disgust, fear, happiness, sadness, and neutral[5]. This emotional speech corpus is probably the most often used database in the context of emotion recognition from speech, and also one of the few for which some results can be compared. We randomly choose 70% utterances of each emotion classification to construct the training dataset, and use the other 30% utterances for test. A summary of the emotion class distribution can be found in Table 2.

### B. Interactive Emotional Dyadic Motion Capture Database

The IEMOCAP database is a well-known dataset for speech emotion recognition comprising of acted and spontaneous multimodal expressive dyadic interactions. The design of the database assumed that by exploiting the context of dyadic interactions between actors, a more natural and richer emotional display would be elicited than in speech read by a single subject. Furthermore, the use of scripted and emotionally targeted improvisational scenarios allowed us to collect a varied and balanced database. The database contains 12 hours of data (audio/video recording and motion capture trajectories of facial markers), split into 5min dyadic interactions between 5 professional female-male actor pairs. Each session consists of a different dyad of male-female actors performing scripted plays and engaging in spontaneous improvised dialogs elicited through affective scenario prompts. Sentences are annotated by at least three Naive humans with categorical emotion labels. In order to match experimental conditions in previously reported categorical emotion recognition studies on USC IEMOCAP [4,9], we consider only 5480 sentences with majority agreement over the emotion classes of: anger, happiness, sadness and neutral. In this paper, We randomly choose 70% utterances of each emotion classification to construct the training dataset,

and use the other 30% utterances for test. The length of an utterance is unequal and its mean is 4s. A summary of the emotion class distribution can be found in Table 2.

## VIII. RESULTS

The Emo-Db database is recorded using the Sennheiser MKH 40 P48 microphone, with the sampling frequency of 16 kHz. Samples are stored as 16 bit numbers. The speech data sampled at 16 kHz from IEMOCAP. The speech data from the two databases is converted into frames with a 18-ms window sliding at 8-ms each time.

First, the signal is passed through a pre-processing system that normalizes amplitude, reduces the amount of noise and extracts combined MFCC parameters [10]. We used the un-weighted average recall (UAR) as a performance metric for emotion recognition results. The UAR is a metric that has been used as the standard measurement in the INTERSPEECH Emotion Challenges. This is the average of the results for each emotion class. With SVMs, normalizing the dynamic ranges of the supervector elements is also crucial since SVMs are not scale invariant.

### A. Results on Emo-DB

Different Gaussian mixtures (4, 8, 16, 32, 64, 128, 256 and 512) are set to compare which Gaussian mixture model can attain a high classification performance. Table 3 shows the accuracy obtained in this experiment, for different number of mixtures in Linear supervector kernel and Non Linear supervector kernel. As expected, the accuracy of the system increases with the number of mixtures, until saturation. The highest recognition is 128 GMM mixtures for linear supervector kernel and is 256 for non linear supervector kernel.

In the next step, we evaluated GUMI kernel and GLDS kernel among all emotions from Emo-Db. From figure 3, we come up with two conclusions: Firstly, it is shown that the GUMI kernel achieves the best accuracy. Secondly, recognition rate of various emotions are different. Angry is the best emotion to be recognized but Neutral speech is the hardest style to be recognized.

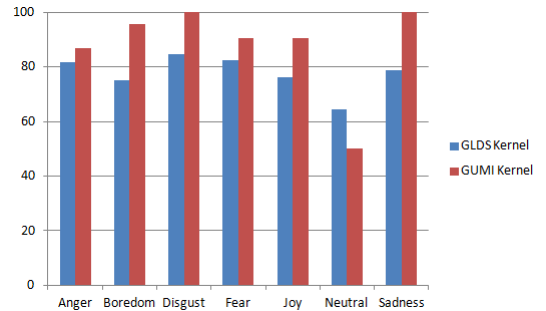


Fig. 3. GLDS Kernel vs GUMI Kernel Performances on Emo-Db.

### B. Results on IEMOCAP

Four confusion matrices are shown in Tables 5-8 for linear GMM kernel, non linear GMM kernel, GUMI kernel and

TABLE I  
EMO-DB DATABASE: NUMBER OF EMOTION UTTERANCES PER CATEGORY.

	Anger	Boredom	Disgust	Fear	Joy	Sadness	Neutral	Total
Total	128	81	44	69	71	45	62	500
Training set	90	57	31	48	50	31	43	350
Testing set	38	24	13	21	21	14	19	150

TABLE II  
IEMOCAP DATABASE: NUMBER OF EMOTION UTTERANCES PER CATEGORY.

	Anger	Joy	Sadness	Neutral	Total
Total	1083	1630	1083	1683	5480
Training set	758	1141	758	1178	3823
Testing set	325	489	325	505	1657

TABLE III  
RECOGNITION RATES FOR LINEAR SUPERVECTOR KERNEL AND NON LINEAR SUPERVECTOR KERNEL USING DIFFERENT NUMBER OF GMM MIXTURES.

UAR(%)	4	8	16	32	64	128	256	512
Linear Kernel	58	61.33	72.67	75.33	82.67	80.67	80	79.7
NON Linear Kernel	54.72	60.89	71.98	74.9	82.97	82	82	80

GLDS kernel, respectively trained with IEMOCAP database. The accuracy column lists per class recognition rates. Firstly, the optimum system is the GUMI kernel (bhattacharya distance), whose recognition rate is better than any other systems. It is also shown that most emotions can be correctly recognized with above 70% accuracy, with the exception of joy, which forms the most notable confusion pair with anger, though they are of opposite valence in the arousalvalence space. This might be due to the fact that arousal is more easily recognized than valence. The most common mistake is between sadness and neutral speech.

## IX. CONCLUSION

Our approaches to robust speech emotion recognition are largely geared to common approaches to automatic speaker recognition. I.e., we concentrate in this paper, our effort to explore the role of GMM supervectors modeling in a discriminative framework. This paper takes a range of dynamic sequence kernels and compares their performance over a range of different emotions. An extensive set of experiments were conducted using two kinds of databases (EMO-Db and IEMOCAP). We find that for Emo-Db, linear GMM kernel outperforms the other kernels, while, for IEMOCAP, GUMI kernel gives the best accuracies. In the future, we will study the predictive power of other sequence kernel methods (e.g., probabilistic kernel). We will also work on building a robust emotion recognition system that utilize different types of prosodic features (e.g., pitch, loudness, energy). We plan also to examine whether our findings generalize to other emotional databases of Arabic dialogues.

## REFERENCES

- [1] Trabelsi, I., Ben Ayed, D. and Ellouze, N. (2013) 'Improved Frame Level Features and SVM Supervectors Approach for the Recognition of Emotional States from Speech: Application to categorical and dimensional states,' *International Journal of Image, Graphics and Signal Processing (IJIGSP)*, vol.5, pp. 813.
- [2] HERM, Ota, SCHMITT, Alexander, et LISCOMBE, Jackson. When calls go wrong: How to detect problematic calls based on log-files and emotions?. In : Ninth Annual Conference of the International Speech Communication Association. 2008. G Goel, V., Olsen, P. A., Rennie, S. J., Huang, J. (2014). U.S. Patent No. 8,738,376. Washington, DC: U.S. Patent and Trademark Office.
- [3] Vapnik, V. (1995) *The nature of statistical learning theory*. Springer-Verlag, New York.
- [4] Trabelsi I. and Ben Ayed D. (2012) On the use of different feature extraction methods for linear and non linear kernels. In Proc. Of Sciences of Electronics, Technologies of Information and Telecommunications (SETIT 2012), 797-802.
- [5] Burkhardt, F., et al. (2005) 'A Database of German Emotional Speech.' *Proc Interspeech. Lisbon, Portugal*, pp. 1517-20.
- [6] Trabelsi, I., Ayed, D. B. (2013). A Multi Level Data Fusion Approach for Speaker Identification on Telephone Speech. *International Journal of Signal Processing, Image Processing Pattern Recognition*, 6(2).
- [7] BURGESS, J.C. (1998) 'A tutorial on support vector machines for pattern recognition.' in *Data mining and knowledge discovery*, vol. 2, no 2, pp. 121-167.
- [8] CHANG, C-C. and LIN, C.J. (2011) 'LIBSVM: a library for support vector machines.' *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 2, no 3, p. 27.
- [9] Busso, C., Bulut, M., Lee, C.-C., Kazemzadeh, A., Mower, E., Kim, S., Chang, J., Lee, S., Narayanan, S., 2008. IEMOCAP: Interactive emotional dyadic motion capture database
- [10] Cullen, C., Vaughan, B., Kousidis, S. (2008, January). Emotional speech corpus construction, annotation and distribution. In *Conference papers* (p. 20).
- [11] Malekzadeh, M., Mustafa, M. B., Lahsasna, A. (2015). A review of emotion regulation in intelligent tutoring systems. *Educational Technology Society*, 18(4), 435-445.
- [12] Black, M. P., Katsamanis, A., Baucom, B. R., Lee, C. C., Lammert, A. C., Christensen, A., ... Narayanan, S. S. (2013). Toward automating a human behavioral coding system for married couples interactions using speech acoustic features. *Speech Communication*, 55(1), 1-21.
- [13] Agapi, C., Mandalia, B. D., Mansey, P. P. (2014). U.S. Patent No. 8,654,937. Washington, DC: U.S. Patent and Trademark Office.
- [14] EYBEN, Florian, WLLMER, Martin, POITSCHKE, Tony, et al. Emotion on the roadnecessity, acceptance, and feasibility of affective computing in the car. *Advances in human-computer interaction*, 2010, vol. 2010.
- [15] Park, C. H., Park, H. (2008). A comparison of generalized linear discriminant analysis algorithms. *Pattern Recognition*, 41(3), 1083-1097.
- [16] Trinh, T. D., Bui, N. N., Min, S. H., Kim, J. Y. (2013). Audio event

TABLE IV  
CONFUSION MATRIX OF THE EMOTION RECOGNITION SYSTEM: LINEAR GMM KERNEL .

	Anger	Joy	Sadness	Neutral	Accuracy
Anger	270	35	5	15	83.07
Joy	100	301	36	62	61.55
Sadness	20	5	250	50	76.92
Neutral	27	28	150	300	59.40

TABLE V  
CONFUSION MATRIX OF THE EMOTION RECOGNITION SYSTEM: NON LINEAR GMM KERNEL .

	Anger	Joy	Sadness	Neutral	Accuracy
Anger	269	34	6	15	82.76
Joy	100	305	19	65	62.37
Sadness	10	7	256	52	78.76
Neutral	30	21	100	254	70.09

TABLE VI  
CONFUSION MATRIX OF THE EMOTION RECOGNITION SYSTEM: GUMI KERNEL .

	Anger	Joy	Sadness	Neutral	Accuracy
Anger	270	33	6	16	83.07
Joy	106	285	28	70	58.28
Sadness	12	6	252	55	77.53
Neutral	10	14	111	370	73.26

TABLE VII  
CONFUSION MATRIX OF THE EMOTION RECOGNITION SYSTEM: GLDS KERNEL .

	Anger	Joy	Sadness	Neutral	Accuracy
Anger	271	29	6	19	83.38
Joy	168	284	20	12	58.07
Sadness	13	12	245	52	75.38
Neutral	36	53	172	244	48.31

classification using SVM with GMM-UBM supervectors., 11(11), 91-98.

- [17] Bui, N. N., Kim, J. Y., Trinh, T. D. (2014). A Non-linear GMM KL and GUMI Kernel for SVM Using GMM-UBM Supervector in Home Acoustic Event Classification. IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences, 97(8), 1791-1794.