



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Emotion Recognition in Low-Resource Settings

Citation for published version:

Haider, F, Pollak, S, Albert, P & Luz, S 2021, 'Emotion Recognition in Low-Resource Settings: An Evaluation of Automatic Feature Selection Methods', *Computer Speech and Language*, vol. 65, 101119. <https://doi.org/10.1016/j.csl.2020.101119>

Digital Object Identifier (DOI):

[10.1016/j.csl.2020.101119](https://doi.org/10.1016/j.csl.2020.101119)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Computer Speech and Language

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Emotion Recognition in Low-Resource Settings: An Evaluation of Automatic Feature Selection Methods

Fasih Haider

Usher Institute, Edinburgh Medical School, the University of Edinburgh, UK

Senja Pollak

*Jozef Stefan Institute, Ljubljana, Slovenia and Usher Institute, Edinburgh Medical School,
the University of Edinburgh, UK*

Pierre Albert

Usher Institute, Edinburgh Medical School, the University of Edinburgh, UK

Saturnino Luz

Usher Institute, Edinburgh Medical School, the University of Edinburgh, UK

Abstract

Research in automatic affect recognition has seldom addressed the issue of computational resource utilization. With the advent of ambient intelligence technology which employs a variety of low-power, resource-constrained devices, this issue is increasingly gaining interest. This is especially the case in the context of health and elderly care technologies, where interventions may rely on monitoring of emotional status to provide support or alert carers as appropriate. This paper focuses on emotion recognition from speech data, in settings where it is desirable to minimize memory and computational requirements. Reducing the number of features for inductive inference is a route towards this goal. In this study, we evaluate three different state-of-the-art feature selection methods: Infinite Latent Feature Selection (ILFS), ReliefF and Fisher (generalized Fisher score), and compare them to our recently proposed feature selection method named ‘Active Feature Selection’ (AFS). The evaluation is performed on three emotion recognition data sets (EmoDB, SAVEE and EMOVO) using two standard acoustic paralinguistic feature sets (i.e. eGeMAPs and emobase). The

results show that similar or better accuracy can be achieved using subsets of features substantially smaller than the entire feature set. A machine learning model trained on a smaller feature set will reduce the memory and computational resources of an emotion recognition system which can result in lowering the barriers for use of health monitoring technology.

Keywords: Feature Engineering, Feature Selection, Emotion Recognition, Affective Computing, Prosodic analysis, Cognitive Health Monitoring

1. Introduction

Speech signals are used in a number of automatic prediction tasks, including cognitive state detection [1], cognitive load estimation [2], presentation quality assessment [3] and emotion recognition [4, 5]. Emotional/affective states could have influence on health and intervention outcomes. Positive emotions have been linked with health improvement, while negative emotions may have negative impact [6]. For example, long term bouts of negative emotions are predisposing factors for depression (ibid.), while positive emotions-related humour and optimism have been linked with positive effects on the immune system and cardiovascular health [7]. Emotion recognition has been used in applications in the domain of health technologies, including mental health assessment and beyond [8, 9, 10].

Applications using speech usually extract emotions as an additional signal in complex systems, such as in ambient intelligence (AmI) [11], depression recognition [9], and longitudinal cognitive status assessment [12]. These approaches employ very high-dimensional feature spaces consisting of large numbers of potentially relevant acoustic features, usually obtained by applying statistical functionals to basic, energy, spectral and voicing related acoustic descriptors [13] extracted from speech intervals lasting a few seconds [14]. Although there is no general consensus on what the ideal set of features should be, this “brute-force” approach of employing as many features as possible seems to outperform alternative (Markovian) approaches to modelling temporal dynamics on the classifier

level [15]. However, the use of such high-dimensional data sets poses challenges for prediction, as they suffer from the so-called “curse of dimensionality”, high degree of redundancy in the feature set, and a large number of features with poor descriptive value. Su and Luz, for instance, noted that in a cognitive load prediction data set about 4% of a feature set of over 250 features had a standard deviation of less than 0.01 and therefore contributed negligibly to the classification task [12]. Moreover, processing of very large numbers of features presents computational challenges for the low-power, low-cost devices such as the Raspberry Pi Zero ¹, which are often used in AmI applications.

The main contribution of this study is the evaluation of different state of the art feature selection methods, including our Active Feature Selection (AFS) method, on the emotion recognition from speech, which has, to the best of our knowledge, not yet been systematically explored. This study extends our previous work [16], where we first introduced the novel AFS method and tested it on the ICMI Challenge on Eating Conditions Recognition [17].

2. Background and Related Work

The automatic identification of emotions in speech is a challenging task, and identifying relevant acoustic features and systematic comparative evaluations has proved difficult [18]. In 2016, the eGeMAPs set [19] (see Section 4.2) was designed based on features’ potential to reflect affective processes and their theoretical significance. It was proposed to set a common ground of emotion-related speech features, which has become since then a *de-facto* standard. The set of target emotions has mostly been fixed around the ‘Big Six’, and similarly, evaluations are more and more frequently performed on a number of publicly available corpora (see Section 4.1). In the health domain, feature selection methods for speech processing have been applied to determine the most discriminant features in support of automatization efforts, as for instance in the

¹<https://www.raspberrypi.org/products/raspberry-pi-zero/> (last accessed January 2019)

50 assessment of patients with pre-dementia and Alzheimer’s disease [20, 21] or for
the detection of sleep apnea [22]. The automatic emotion recognition task has
gained attention in the past few years [23, 24]. This task has been addressed
through processing of facial, speech, body movements and biometric informa-
tion [25, 26, 27, 28, 1]. Numerous studies [25, 26, 28, 29, 30, 31] extract audio
55 features with OpenSMILE using *de-facto* standard presets: IS10, GeMAPS,
eGeMAPs, Emobase.

The reviewed literature suggests that although the accuracy of various ma-
chine learning approaches in this area is promising, automatic dimensionality
reduction has focused largely on the removal of noisy or redundant features,
60 with less attention paid to computational resource utilisation [1, 25, 26, 27, 28,
29, 30, 31].

There are many dimensionality reduction methods: some are feature selec-
tion methods which require labelled data, and some are feature transformation
methods which do not require labelled data. The former includes methods such
65 as correlation based feature selection and Fisher feature selection [32, 33], while
the latter includes, for instance, principal component analysis (PCA), inde-
pendent component analysis (ICA) [34] among others. Recently, efforts have
focused on reducing dimensionality using PCA to improve the results for emo-
tion recognition from speech [35, 36, 37, 38] in different settings such as noisy
70 setting [36]. Dimensionality reduction using feature selection methods, on the
other hand, are less explored in this area.

3. Feature Selection Methods

In this section we will briefly describe the feature selection methods used in
this study along with our AFS method. We have selected three state of the art
75 feature selection methods. The motivation behind using these methods here is
their robust performance in a number of tasks [39].

3.1. Infinite Latent Feature Selection (ILFS)

The ILFS method [39] performs cross-validation on an unsupervised ranking of features. At a pre-processing stage, each feature is represented by a descriptor reflecting how discriminative it is. A probabilistic latent graph containing each
80 feature is built. Weighted edges model pairwise relations among feature distributions, created using probabilistic latent semantic analysis. The relevance of each feature is computed by looking on its weight in arbitrary set of cues. Each path in the graph represents a selection of features. The final ranking of each
85 feature looks at its redundancy in all the possible feature subsets, selecting the most discriminative and relevant features. The evaluation on a range of different tasks (e.g. object recognition classification and DNA microarray analysis) confirms its robustness, outperforming other methods on robustness and ranking quality [39].

90 3.2. ReliefF

The ReliefF algorithm [40] which is an adaptation of the Relief feature selection method [41], performs ranking and selection of top scoring features based on their processed score. The score is calculated by weighting features on a random sample of instances. For each instance, the weight vector represents
95 the relevance of each feature amongst the class labels: neighbours are selected from the same class (nearest hits) and from each different class (nearest misses). The weight of each feature increases when the difference with its nearest hits is low and with its nearest misses is high. Each weight vector is combined in a global relevance vector. The final subset is constituted of all the features with
100 relevance greater than a manually set threshold. ReliefF is a common method of Feature Selection which has been continuously improved since its first publication [41, 42].

3.3. Generalized Fisher score (Fisher)

The generalized Fisher score [33] is a generalization of the Fisher score to
105 take into account redundancy and combination of features. A subset of features

is sought which maximizes the lower bound of the traditional Fisher score. A combination of features is evaluated, and redundant features discarded. A quadratically constrained linear programming (QCLP) is solved with a cutting plane algorithm. At each iteration, a multiple kernel learning is solved by a
110 multivariate ridge regression followed by a projected gradient descent to update the kernel weights. The method produces state of the art results, outperforming many feature selection methods while having a lower complexity [33].

3.4. Active feature selection method

An Active Feature Selection method, which divides a feature set into subsets,
115 has been recently introduced [16]. The term ‘active’ is used because compared to other approaches it evaluates feature subsets and not each feature separately, so that different features actively contribute to the feature selection. While clustering is employed, AFS does not cluster instances but dimensions. Our hypothesis is that noisy features have common characteristics that differ from
120 those of informative features, and that clustering will divide the features into subsets according to such common characteristics. This involves clustering the data set into N clusters (where $N = 5, 10, 15, \dots, 100$) using self-organizing maps (SOM) with 200 iterations and batch training [43], and then evaluating the discrimination power of the features from each cluster C_N using leave one subject
125 out (LOSO) cross-validation, as shown in Figure 1. The cluster with the highest validation accuracy is selected (see Figure 6 in Section 5).

4. Experimentation

The section describes the datasets and their characteristics along with acoustic feature extraction and classification methods.

4.1. Data sets

130 Three corpora were selected for their shared characteristics and public availability: EmoDB, SAVEE, and EMOVO. They consist of recorded acted performances, annotated using the well-known and widely used *Big Six* set of anno-

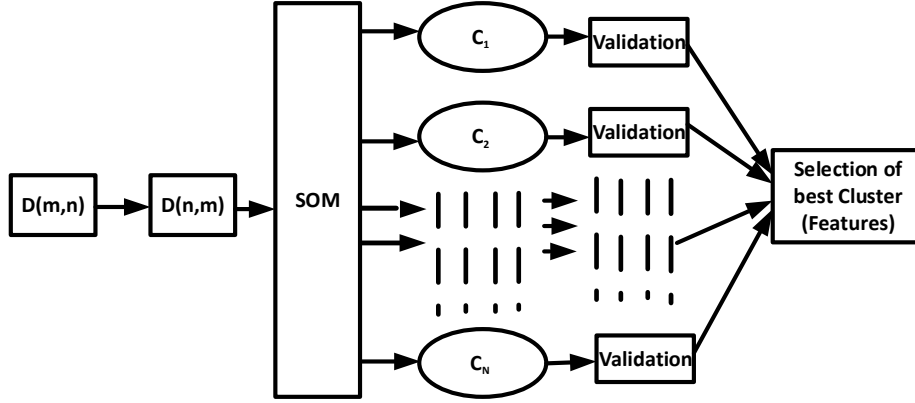


Figure 1: Active feature selection method: $D(m,n)$ represents the data where m is the total number of training instances and n is the total number of dimensions (988 for *emobase* and 88 for *eGeMAPs*) [16].

tations : anger, disgust, fear, happiness, sadness, surprise + neutral, except in
 135 the older EmoDB data set where boredom was used instead of surprise. Their
 characteristics are summarised in Tables 1 and 2.

Berlin Database of Emotional Speech (EmoDB)

The EmoDB corpus [44] is a data set commonly used in the automatic emotion
 recognition literature. It features 535 acted emotions in German, based on ut-
 140 terances carrying no emotional bias. The corpus was recorded in a controlled
 environment resulting in high quality recordings. Actors were allowed to move
 freely around the microphones, which affected absolute signal intensity. In ad-
 dition to the emotion, each recording was labelled with phonetic transcription
 using the SAMPA phonetic alphabet, emotional characteristics of voice, seg-
 145 mentation of the syllables, and stress. The quality of the data set was evaluated
 by perception tests carried out by 20 human participants. In a first recognition
 test, subjects listened to a recording once before assigning one of the available
 categories, achieving an average recognition rate of 86%. A second naturalness
 test was performed. Documents achieving a recognition rate lower than 80%
 150 or a naturalness rate lower than 60% were discarded from the main corpus,
 reducing the corpus to 535 recordings from the original 800.

Surrey Audio-Visual Expressed Emotion (SAVEE)

SAVEE [45] is an audio-visual data set that was recorded to support the development of an automatic emotion recognition system. The corpus is a set of
155 480 British English utterances. Each actor was recorded for 15 utterances per emotion (3 common utterances recorded for each of the 7 emotions, 2 emotion specific, and 10 generic sentences different for each emotion) and 30 neutral recordings (the 3 common and every emotion specific sentences). No limitation
160 the description of the data set. A qualitative evaluation of the database was run as a perception tests by 10 human subjects. The mean classification accuracy for the audio modality was 66.5%, 88% for the visual modality, and 91.8% for the combined audio-visual modalities.

Italian Emotional Speech Database (EMOVO)

165 The EMOVO corpus [46] is a speech data set featuring recorded emotions from acted performances by 6 persons. Actors were allowed to move freely around the microphones and the volume was manually adjusted, affecting absolute signal intensity. A qualitative evaluation was performed using a discrimination test. Two phrases were selected and, for each, 12 subjects had to choose between two
170 proposed emotions. The mean accuracy for the test was about 80%.

4.2. Volume normalization and feature extraction

We have normalized all the speech utterances' volume into the range [-1:+1] dBFS before any acoustic feature extraction. The motivation for this is to improve the model's robustness against different recording conditions such as
175 distance between microphone and subject. We use the openSMILE [47] toolkit for the extraction of two acoustic feature sets which are widely used for emotion recognition. These are:

emobase: this acoustic feature set contains the MFCC, voice quality, fundamental frequency (F0), F0 envelope, LSP and intensity features along with
180 their first and second order derivatives. In addition, many statistical functions are applied to these features, resulting in a total of 988 features for every speech

Table 1: Main characteristics of the data sets.

Corpus	Size (utterances)	Population	Participants	Emotion categories
EmoDB	535	10 (5 males, 5 females)	German native speakers actors	anger, disgust, fear, joy, sadness, <i>boredom</i> + neutral
SAVEE	480	4 (males)	English native speakers actors	anger, disgust, fear, happiness, sadness, <i>surprise</i> + neutral
EMOVO	588	6 (3 males, 3 females)	Italian native speakers actors	anger, disgust, fear, happiness, sadness, <i>surprise</i> + neutral

Table 2: Distribution of recordings across emotion categories.

Corpus	Neutral	Anger	Disgust	Fear	Happiness	Sadness	Surprise	Boredom
EmoDB	79	127	46	69	71	62	-	81
SAVEE	120	60	60	60	60	60	60	-
EMOVO	84	84	84	84	84	84	84	-

utterance.

eGeMAPs: this feature set contains the F0 semitone, loudness, spectral flux, MFCC, jitter, shimmer, F1, F2, F3, alpha ratio, hammarberg index and
185 slope V0 features including many statistical functions applied to these features, which result in a total of 88 features for every speech utterance [19].

4.3. Classification Method

Classification is performed using Support Vector Machines (SVM) with a linear kernel, SMO solver and cost parameter (box constraint) set to 0.75. This

190 classifier is employed in MATLAB² using the statistics and machine learning toolbox. The feature selection methods are evaluated through LOSO cross-validation, and unweighted average recall (UAR) results are computed.

4.4. Evaluation Criterion

All of the emotion recognition data sets are labeled for seven classes and
195 we have evaluated the classifier using UAR, which corresponds to the average accuracy of all classes. The UAR measure is selected because the datasets are not balanced for emotions. The method with the highest UAR is considered the best. The blind/majority guess for this task results in a 14.3% UAR. As our focus is on feature selection methods, we set the baseline as UAR obtained
200 using the entire feature set.

5. Results and discussion

We have evaluated the three different automatic feature selection methods (ILFS, ReliefF and Fisher) along with our AFS method using two different acoustic feature sets extracted from three different data sets. The results of
205 three feature selection methods are shown in Figure 2. The AFS results are not plotted there, as the AFS does not operate on features iteratively, but on subsets of features determined through SOM. It can be observed that around 30 out of 88 eGeMAPs features and around 100 out of 988 emobase features are sufficient to provide almost the same UAR as the highest achieved UAR for the
210 three data sets. The best results of each feature selection method are shown in Table 3.

The results confirm that a higher accuracy can be achieved using a subset of the feature set than when using the full feature set. The results for each data set can be summarised as follows:

- 215 1. EmoDB: the ILFS method provides better UAR (69.7% for eGeMAPs and 76.9% for emobase) results than the other methods and is able to

²<http://uk.mathworks.com/products/matlab/> (Last accessed: January 2019)

Table 3: Best Unweighted Average Recall (UAR (%)) of feature selection methods and number of selected features (numFeat) are reported. The best UAR (%) results for each feature set are given in bold. The unweighted arithmetic average for each feature selection method is also reported in ‘Average’ column.

Data Set	EmoDB				EMOVO				SAVEE				
Feature Set	eGeMAPs		emobase		eGeMAPs		emobase		eGeMAPs		emobase		
	numFeat	UAR (%)	numFeat	UAR (%)	numFeat	UAR	numFeat	UAR(%)	numFeat	UAR (%)	numFeat	UAR (%)	Mean
Baseline	88	68.5	988	74.6	88	37.4	988	34.4	88	40.8	988	38.1	49.0
ILFS	74	69.7	685	76.9	28	38.1	113	34.7	86	42.0	574	38.8	46.9
reliefF	88	68.5	666	75.3	20	37.8	348	37.1	82	41.4	72	39.3	49.9
Fisher	88	68.5	975	75.2	25	41.0	464	36.2	34	42.4	158	42.4	51.0
AFS	81	68.5	696	75.8	2	39.0	56	36.4	68	40.5	21	37.5	49.6

reduce the number of features (74 out of 88 for eGeMAPs and 685 out of 988 for emobase). The confusion matrix of the best UAR (76.9%) is shown in Figure 3. For eGeMAPs, the AFS method provides an UAR of 68.5% (around 1% lower than ILFS) using 81 features. For emobase, AFS method provides an UAR of 75.8% (around 1% lower than ILFS) using 696 features. With a subset of the eGeMAPs feature set, the reliefF and Fisher methods are not able to improve over the baseline in terms of UAR. However, Figure 2 shows that reliefF and Fisher achieved almost the same UAR as compared to baseline with only 35 eGeMAPs features instead of 88 eGeMAPs features. Hence around 60% reduction in number of features is observed.

2. EMOVO: the Fisher method yields the best UAR (41.0%) using only 25 out of 88 eGeMAPs features, while ReliefF method yields the best UAR (37.1%) for emobase (selecting 348 out of 988 features). The confusion matrix of the best UAR (41.0%) is shown in Figure 4. The results for AFS are slightly lower than the best method (around 2%), but the number of features are significantly lower, compared to other methods. AFS selects only 2 eGeMAPs features out of 88, and 56 emobase feature out of 988, while still reaching an UAR of 39.0% and 36.4%, respectively.
3. SAVEE: the Fisher method again yields the best UAR for eGeMAPs (34 features, and UAR of 42.4%) and emobase (158 features and UAR of

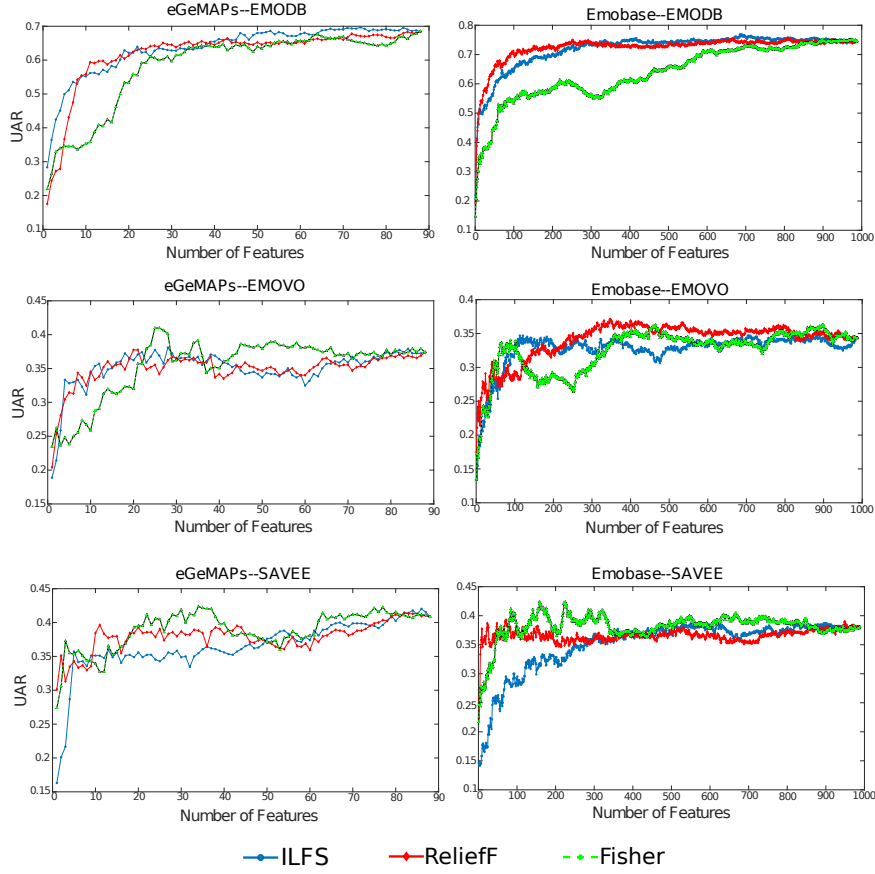


Figure 2: Feature selection methods (ILFS, ReliefF and Fisher) results for all three data sets (EMODB, EMOVO and SAVEE) using two feature sets (eGeMAPs and emobase). Where x-axis represents the number of features and y-axis represents the UAR.

42.4%). The confusion matrix of the best result ($UAR = 42.4\%$) using eGeMAPs features is shown in Figure 5. For eGeMAPs, the results of AFS are slightly lower than the best method (around 2%). For emobase, AFS method yields and UAR of 37.5% (around 5% lower than Fisher) using 21 features.

The machine learning models trained using EmoDB ($UAR=76.9\%$) data provide better UAR than EMOVO (41.0%) and SAVEE (42.9%). This could be due to very high quality nature of the EmoDB data set. The EmoDB data

True Class	Anger	112	0	1	0	14	0	0	Recall
	Boredom	1	67	4	0	0	5	4	88.2%
	Disgust	2	2	36	1	1	3	1	82.7%
	Fear	9	0	0	49	3	4	4	78.3%
	Happiness	25	0	1	6	39	0	0	71.0%
	Neutral	0	7	2	2	3	65	0	54.9%
	Sadness	0	6	4	1	0	1	50	82.3%
	Precision	75.2%	81.7%	75.0%	83.1%	65.0%	83.3%	84.7%	
<div> <div>Anger</div> <div>Boredom</div> <div>Disgust</div> <div>Fear</div> <div>Happiness</div> <div>Neutral</div> <div>Sadness</div> </div>									
Predicted Class									UAR = 76.9 %

Figure 3: Confusion matrix of ILFS Feature selection method for EmoDB data set using emobase feature set.

True Class	Anger	52	2	8	12	8	1	1	Recall
	Disgust	3	19	11	19	13	11	8	61.9%
	Fear	5	4	31	17	9	15	3	22.6%
	Joy	12	3	12	31	9	3	14	36.9%
	Neutral	8	2	7	21	32	12	2	36.9%
	Sadness	2	3	11	4	12	41	11	38.1%
	Surprise	3	5	10	20	3	8	35	48.8%
	Precision	61.2%	50.0%	34.4%	25.0%	37.2%	45.1%	47.3%	
<div> <div>Anger</div> <div>Disgust</div> <div>Fear</div> <div>Joy</div> <div>Neutral</div> <div>Sadness</div> <div>Surprise</div> </div>									
Predicted Class									UAR = 41.0 %

Figure 4: Confusion matrix of Fisher feature selection method for EMOVO data set using eGeMAPs feature set.

set quality was evaluated by 20 human coders with an average recognition rate of 86%, and audio recordings with the inter-coder agreement below 80% were removed (no such measure was taken for EMOVO and SAVEE).

For EMOVO, while the reported accuracy for the test set is 80% (see Section 4.1), one should note that rather than evaluating the full EMOVO data set only two phrases were selected and each coder had to choose only between

True Class	Recall						
	Anger	Disgust	Fear	Happiness	Neutral	Sadness	Surprise
Anger	34	5	2	2	12	3	2
Disgust	13	9	0	1	19	10	8
Fear	7	12	15	9	2	2	13
Happiness	5	14	11	21	3	1	5
Neutral	28	3	0	0	86	3	0
Sadness	10	4	1	0	18	27	0
Surprise	0	12	11	6	2	0	29
Precision							
	35.1%	15.3%	37.5%	53.8%	60.6%	58.7%	50.9%
Predicted Class							
	Anger	Disgust	Fear	Happiness	Neutral	Sadness	Surprise
							UAR = 42.4 %

Figure 5: Confusion matrix of Fisher feature selection method for SAVEE data set using eGeMAPs feature set.

two proposed emotions rather than seven. The fact that our machine learning approach to EMOVO classification is of a seven-class problem explains the much lower results obtained in comparison to human performance.

255 For SAVEE, 10 human subjects evaluated the data set and came up with an accuracy of 66.5% for audio. Our machine learning based models provide promising results as compared to humans subjects. Although they are less accurate than human annotators, we use only acoustic information to automate the process of emotion recognition, while human annotators used both acoustic
260 and linguistic information (i.e. the spoken content).

As shown in Table 3, Generalized Fisher score provides better results in 3 out of 6 cases, ILFS provides better results in 2 out of 6 cases and reliefF provides better results in 1 out of 6 cases, indicating that overall Fisher feature selection provides the best results for the emotion recognition task.

265 The AFS method comes second in 3 out of 6 cases as shown in Table 3. It is also observed that the AFS method provides almost the same results in terms of UAR as the other state of the art feature selection methods, with smaller numbers of dimensions on average. We have note that for the SAVEE data set only 2 out of 88 eGeMAPs features (selected by AFS) provide better

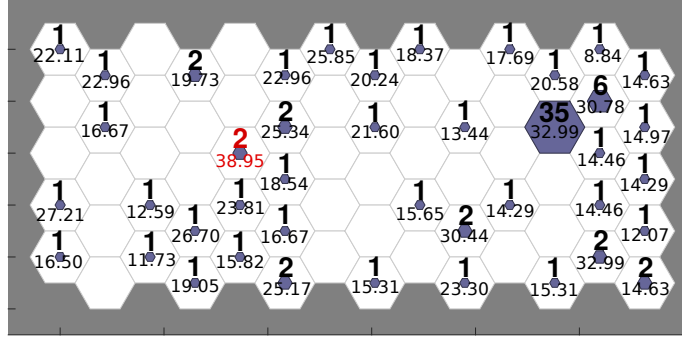


Figure 6: A visualization of AFS method results: number of features present in each cluster (i.e. hexagon or neuron) along with the UAR (%) obtained using eGeMAPs feature set for EMOVO data set. Note that 2 out of 88 features provide better results than other feature subsets.

270 results than reliefF, ILFS and the baseline (i.e. entire feature set). For further insight into these results, we show the evaluation of clusters (feature subsets) using AFS in Figure 6. In this figure we see that there are many clusters which provide better results than the blind guess (14.3%), while the feature cluster selected by AFS contains only 2 features (*hammarbergIndexV_sma3nz_amean* and *hammarbergIndexV_sma3nz_amean*) and leads to the 39.0% UAR. One of the possible lines of future work is to combine features from different clusters to see if this leads to improvement in classification. The AFS method was also evaluated with different numbers of clusters. The best UAR is obtained using 70 clusters for EMOVO dataset. The UAR values for these 70 clusters with their respective numbers of features are shown in Figure 6.

To further evaluate the feature selection methods, we have combined all three data sets which results in a 8-class problem i.e. to recognise (7+1) emotions. The results of this experimentation in LOSO cross-validation setting is shown in Table 4. We have noted that the reliefF method provides the best results for eGeMAPs (46.6%) and emobase (48.0%) feature sets. All three data sets belong to different languages and have different qualities of annotation. Hence, the reliefF method could be a better choice than other methods where the quality and language of data sets are different.

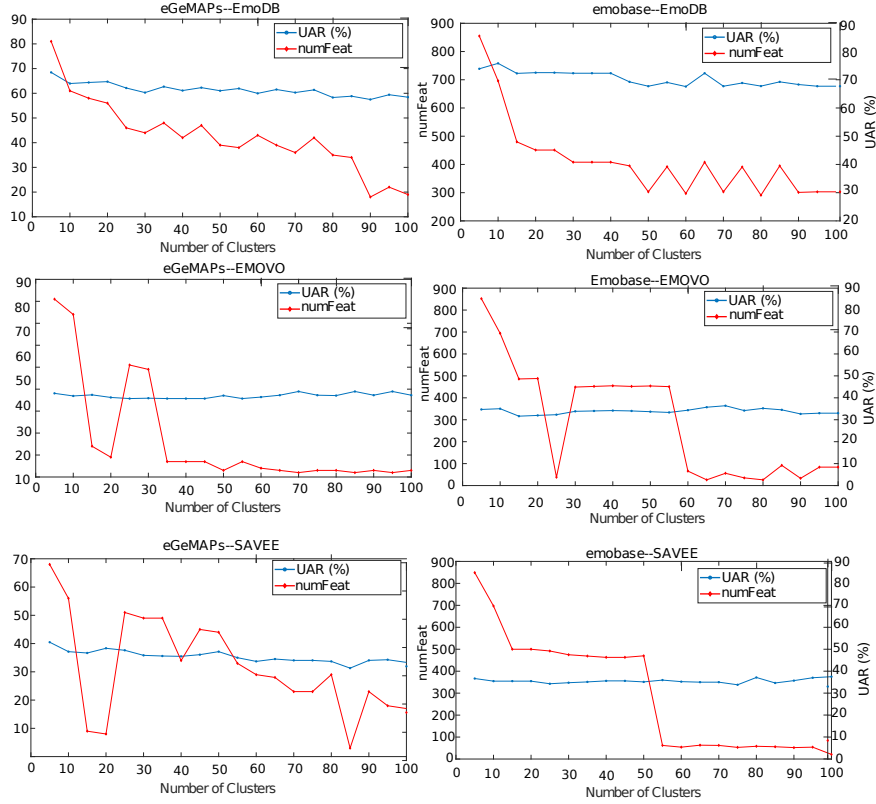


Figure 7: AFS method results: The x-axis represents the number of cluster ($N = 5, 10, 15, \dots, 100$). The y-axis represents the number of features (numFeat) and Unweighted Average Recall (UAR) in % of the best cluster.

In a previous study [16], we demonstrated that the AFS method is able to select a feature subset which provides better results than the entire feature set and the PCA feature set for eating condition recognition. However the results have not been demonstrated in detail as in this study, and the AFS method has not been evaluated on multiple data sets and compared against other feature selection methods to the same extent as in this paper. The present study is therefore a step towards in demonstrating the generalisability of the AFS method. The contribution of this study is not only the evaluation of performance of different feature selection methods but also the assessment of the extent to which AFS, reliefF, Fisher and ILFS can reduce the feature set and therefore

Table 4: Evaluation of feature selection methods for 7+1 emotion recognition task by combining all three data sets. Best Unweighted Average Recall (UAR (%)) and number of selected features (numFeat) are reported. The bold figures indicate the best UAR (%) for each feature set (i.e. eGeMAPs and emobase).

Method	eGeMAPs		emobase	
	numFeat	UAR (%)	numFeat	UAR (%)
Baseline	88	44.4	988	47.4
ILFS	78	45.6	709	47.9
reliefF	44	46.6	732	48.0
Fisher	53	45.3	822	47.9
AFS	79	43.8	835	47.2

select small enough subsets which will impose lower computational demands
 300 on low resource systems, while preserving or improving emotion recognition
 performance, in comparison to full feature sets.

6. Conclusion

This study evaluated three state-of-the-art feature selection methods, namely,
 ILFS, reliefF and generalized Fisher score for emotion recognition, along with
 305 the recently proposed AFS method. It employed three different emotion recog-
 nition data sets from three different languages. The results show that higher
 UAR can be achieved using reduced feature sets. Generally, around 30 out of
 88 eGeMAPs and 100 out of 988 emobase features are sufficient to obtain al-
 most the same UAR as a full feature set. The Fisher feature selection method
 310 provided the best averaged UAR across all three data sets (51.0%) and two
 feature sets compared to the 49.0% averaged UAR for the full feature set base-
 line. However the reliefF method outperformed the other methods when all the
 data sets were combined. These findings are relevant to the development of
 machine learning models for machines with low computational resources. The
 315 AFS method provides competitive results in relation to the state of the art in
 feature select. AFS currently uses only features present in one cluster. For
 future studies, we will explore methods to rank the clusters of features and do

fusion of different clusters for possible accuracy improvements. Other possible
avenues for future work include testing the AFS on other modalities in addition
320 to speech.

7. Acknowledgement

This research is funded by the European Union’s Horizon 2020 research
program, under grant agreement No 769661, towards the SAAM project. PA is
supported by the Medical Research Council (MRC). The work of S. Pollak was
325 partially supported by the Slovenian Research Agency (ARRS) core research
programme *Knowledge Technologies* (P2-0103).

References

- [1] H. Akira, F. Haider, L. Cerrato, N. Campbell, S. Luz, Detection of cognitive
states and their correlation to speech recognition performance in speech-
330 to-speech machine translation systems, in: Proceedings of the 16th An-
nual Conference of the International Speech Communication Association,
INTERSPEECH 2015, International Speech Communications Association,
2015, pp. 2539–2543.
- [2] B. Schuller, S. Steidl, A. Batliner, J. Epps, F. Eyben, F. Ringeval,
335 E. Marchi, Y. Zhang, The INTERSPEECH 2014 computational paralin-
guistics challenge: Cognitive & physical load, in: Proceedings of the 15th
Annual conference of the International Speech Communication Association,
INTERSPEECH 2014, International Speech Communications Association,
2014, pp. 427–431.
- [3] F. Haider, L. Cerrato, N. Campbell, S. Luz, Presentation quality assessment
340 using acoustic information and hand movements, in: Proceedings of the
IEEE International Conference on Acoustics, Speech and Signal Processing
(ICASSP), Institute of Electrical and Electronics Engineers (IEEE), 2016,
pp. 2812–2816.

- 345 [4] M. El Ayadi, M. S. Kamel, F. Karray, Survey on speech emotion recognition: Features, classification schemes, and databases, *Pattern Recognition* 44 (3) (2011) 572–587.
- [5] B. Schuller, A. Batliner, S. Steidl, D. Seppi, Recognising realistic emotions and affect in speech: State of the art and lessons learnt from the first
350 challenge, *Speech Communication* 53 (910) (2011) 1062–1087.
- [6] N. S. Consedine, J. T. Moskowitz, The role of discrete emotions in health outcomes: A critical review, *Applied and Preventive Psychology* 12 (2) (2007) 59–75.
- [7] J. E. Dimsdale, Psychological stress and cardiovascular disease, *Journal of the American College of Cardiology* 51 (13) (2008) 1237–1246.
355
- [8] M. Valstar, B. Schuller, K. Smith, F. Eyben, B. Jiang, S. Bilakhia, S. Schnieder, R. Cowie, M. Pantic, AVEC 2013: the continuous audio/visual emotion and depression recognition challenge, in: *Proceedings of the 3rd ACM international workshop on Audio/visual emotion challenge (AVEC)*, Association for Computing Machinery, 2013, pp. 3–10.
360
- [9] B. Desmet, V. Hoste, Emotion detection in suicide notes, *Expert Systems with Applications* 40 (16) (2013) 6351–6358.
- [10] F. Haider, S. De La Fuente Garcia, P. Albert, S. Luz, Affective speech for alzheimer’s dementia recognition., in: D. Kokkinakis, K. Lundholm Fors, C. Themistocleous, M. Antonsson, M. Eckerström (Eds.), *LREC: Resources and ProcessIng of linguistic, para-linguistic and extra-linguistic Data from people with various forms of cognitive/psychiatric/developmental impairments (RaPID)*, European Language Resources Association (ELRA), 2020, pp. 67–73.
365
- 370 [11] L. Y. Mano, B. S. Faial, L. H. Nakamura, P. H. Gomes, G. L. Libralon, R. I. Meneguete, P. R. Geraldo Filho, G. T. Giancristofaro, G. Pessin, B. Krishnamachari, Exploiting IoT technologies for enhancing Health Smart Homes

through patient identification and emotion recognition, *Computer Communications* 89 (2016) 178–190.

- 375 [12] J. Su, S. Luz, Predicting cognitive load levels from speech data, in: *Recent Advances in Nonlinear Speech Processing*, Springer, 2016, pp. 255–263.
- [13] F. Eyben, M. Wöllmer, B. Schuller, Openearintroducing the munich open-source emotion and affect recognition toolkit, in: *Proceedings of the 3rd International Conference on Affective Computing and Intelligent Interaction and Workshops (ACII)*, IEEE, 2009, pp. 1–6.
- 380 [14] D. Ververidis, C. Kotropoulos, Emotional speech recognition: Resources, features, and methods, *Speech Communication* 48 (9) (2006) 1162–1181.
- [15] F. Weninger, F. Eyben, B. W. Schuller, M. Mortillaro, K. R. Scherer, On the Acoustics of Emotion in Audio: What Speech, Music, and Sound have in Common, *Frontiers in Psychology* 4.
- 385 [16] F. Haider, S. Pollak, E. Zarogianni, S. Luz, SAAMEAT: active feature transformation and selection methods for the recognition of user eating conditions, in: *Proceedings of the 2018 International Conference on Multimodal Interaction (ICMI)*, ACM, Association for Computing Machinery, 2018, pp. 564–568.
- 390 [17] B. Schuller, S. Steidl, A. Batliner, S. Hantke, F. Hnig, J. R. Orozco-Arroyave, E. Nth, Y. Zhang, F. Weninger, The INTERSPEECH 2015 computational paralinguistics challenge: nativeness, parkinson’s & eating condition, in: *Proceedings of the 16th Annual Conference of the International Speech Communication Association, INTERSPEECH 2015*, 2015, pp. 478–482.
- 395 [18] C.-N. Anagnostopoulos, T. Iliou, I. Giannoukos, Features and classifiers for emotion recognition from speech: a survey from 2000 to 2011, *Artificial Intelligence Review* 43 (2) (2015) 155–177.

- 400 [19] F. Eyben, K. R. Scherer, B. W. Schuller, J. Sundberg, E. André, C. Busso, L. Y. Devillers, J. Epps, P. Laukka, S. S. Narayanan, et al., The geneva minimalistic acoustic parameter set (gemaps) for voice research and affective computing, *IEEE Transactions on Affective Computing* 7 (2) (2016) 190–202.
- 405 [20] A. Knig, A. Satt, A. Sorin, R. Hoory, O. Toledo-Ronen, A. Derreumaux, V. Manera, F. Verhey, P. Aalten, P. H. Robert, Automatic speech analysis for the assessment of patients with predementia and Alzheimer’s disease, *Alzheimer’s & Dementia: Diagnosis, Assessment & Disease Monitoring* 1 (1) (2015) 112–124.
- 410 [21] F. Haider, S. de la Fuente, S. Luz, An assessment of paralinguistic acoustic features for detection of alzheimer’s dementia in spontaneous speech, *IEEE Journal of Selected Topics in Signal Processing* 14 (2) (2020) 272–281.
- [22] E. Goldshtein, A. Tarasiuk, Y. Zigel, Automatic detection of obstructive sleep apnea using speech signals, *IEEE Transactions on biomedical engineering* 58 (5) (2011) 1373–1382.
- 415 [23] A. Dhall, A. Kaur, R. Goecke, T. Gedeon, Emotiw 2018: Audio-video, student engagement and group-level affect prediction, in: *Proceedings of the 2018 on International Conference on Multimodal Interaction*, ACM, 2018, pp. 653–656.
- 420 [24] A. Dhall, R. Goecke, S. Ghosh, J. Joshi, J. Hoey, T. Gedeon, From individual to group-level emotion recognition: Emotiw 5.0, in: *Proceedings of the 19th ACM International Conference on Multimodal Interaction (ICMI)*, ICMI 2017, Association for Computing Machinery, 2017, pp. 524–528.
- 425 [25] B. Knyazev, R. Shvetsov, N. Efremova, A. Kuharenko, Convolutional neural networks pretrained on large face recognition datasets for emotion classification from video, *arXiv preprint arXiv:1711.04598*.

- [26] F. Haider, L. S. Cerrato, S. Luz, N. Campbell, Attitude recognition of video bloggers using audio-visual descriptors, in: Proceedings of the Workshop on Multimodal Analyses Enabling Artificial Agents in Human-Machine Interaction, MA3HMI 2016, Association for Computing Machinery, 2016, pp. 38–42.
- [27] N. A. Madzlan, Y. Huang, N. Campbell, Automatic classification and prediction of attitudes: Audio-visual analysis of video blogs, in: Proceedings of the International Conference on Speech and Computer, Springer, 2015, pp. 96–104.
- [28] P. Hu, D. Cai, S. Wang, A. Yao, Y. Chen, Learning supervised scoring ensemble for emotion recognition in the wild, in: Proceedings of the 19th ACM International Conference on Multimodal Interaction (ICMI), ICMI 2017, Association for Computing Machinery, 2017, pp. 553–560.
- [29] V. Vielzeuf, S. Pateux, F. Jurie, Temporal multimodal fusion for video emotion classification in the wild, in: Proceedings of the 19th ACM International Conference on Multimodal Interaction (ICMI), ICMI 2017, Association for Computing Machinery, 2017, p. 569576.
- [30] S. Wang, W. Wang, J. Zhao, S. Chen, Q. Jin, S. Zhang, Y. Qin, Emotion recognition with multimodal features and temporal models, in: Proceedings of the 19th ACM International Conference on Multimodal Interaction (ICMI), ICMI 2017, Association for Computing Machinery, 2017, pp. 598–602.
- [31] X. Ouyang, S. Kawaai, E. G. H. Goh, S. Shen, W. Ding, H. Ming, D.-Y. Huang, Audio-visual emotion recognition using deep transfer learning and multiple temporal models, in: Proceedings of the 19th ACM International Conference on Multimodal Interaction (ICMI), ACM, Association for Computing Machinery, 2017, pp. 577–582.
- [32] M. A. Hall, Correlation-based feature selection for machine learning, Ph.D. thesis, The University of Waikato (1999).

- [33] Q. Gu, Z. Li, J. Han, Generalized fisher score for feature selection, arXiv preprint arXiv:1202.3725.
- [34] J. Wang, C.-I. Chang, Independent component analysis-based dimensionality reduction with applications in hyperspectral image analysis, *IEEE transactions on geoscience and remote sensing* 44 (6) (2006) 1586–1600.
- [35] N. P. Jagini, R. R. Rao, Exploring emotion specific features for emotion recognition system using pca approach, in: *Proceedings of the International Conference on Intelligent Computing and Control Systems (ICICCS)*, IEEE, 2017, pp. 58–62.
- [36] P. K. Aher, S. D. Daphal, A. N. Cheeran, Analysis of feature extraction techniques for improved emotion recognition in presence of additive noise, in: *Proceedings of the International Conference on Computation System and Information Technology for Sustainable Solutions (CSITSS)*, IEEE, 2016, pp. 350–354.
- [37] S. Wang, X. Ling, F. Zhang, J. Tong, Speech emotion recognition based on principal component analysis and back propagation neural network, in: *Proceedings of the International Conference on Measuring Technology and Mechatronics Automation (ICMTMA)*, Vol. 3, IEEE, 2010, pp. 437–440.
- [38] F. Haider, F. A. Salim, O. Conlan, S. Luz, An active feature transformation method for attitude recognition of video bloggers, in: *Proc. Interspeech 2018*, 2018, pp. 431–435. doi:10.21437/Interspeech.2018-1222. URL <http://dx.doi.org/10.21437/Interspeech.2018-1222>
- [39] G. Roffo, S. Melzi, U. Castellani, A. Vinciarelli, Infinite latent feature selection: A probabilistic latent graph-based ranking approach, in: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, IEEE, 2017, pp. 1407–1415.
- [40] I. Kononenko, E. Šimec, M. Robnik-Šikonja, Overcoming the myopia of

inductive learning algorithms with ReliefF, *Applied Intelligence* 7 (1) (1997) 39–55.

- 485 [41] K. Kira, L. A. Rendell, et al., The feature selection problem: Traditional methods and a new algorithm, in: *Aaai*, Vol. 2, 1992, pp. 129–134.
- [42] M. Robnik-Šikonja, I. Kononenko, An adaptation of Relief for attribute estimation in regression, in: *Proceedings of the Fourteenth International Conference on Machine Learning (ICML)*, Vol. 5 of ICML 1997, 1997, pp. 296–304.
- 490 [43] T. Kohonen, The self-organizing map, *Neurocomputing* 21 (1-3) (1998) 1–6.
- [44] F. Burkhardt, A. Paeschke, M. Rolfes, W. F. Sendlmeier, B. Weiss, A database of german emotional speech, in: *Proceedings of the ninth European Conference on Speech Communication and Technology*, 2005, pp. 1516–1520.
- 495 [45] S. Haq, P. Jackson, Speaker-dependent audio-visual emotion recognition, in: *Proceedings of the International Conference on Auditory-Visual Speech Processing (AVSP)*, 2009, pp. 53–58.
- [46] G. Costantini, I. Iaderola, A. Paoloni, M. Todisco, Emovo corpus: an italian emotional speech database, in: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC)*, LREC 2014, European Language Resources Association (ELRA), 2014, pp. 3501–3504.
- 500 [47] F. Eyben, F. Weninger, F. Groß, B. Schuller, Recent developments in opensmile, the munich open-source multimedia feature extractor, in: *Proceedings of the 21st ACM international conference on Multimedia*, ACM, Association for Computing Machinery, 2013, pp. 835–838.
- 505