

Leveraging Inter-rater Agreement for Audio-Visual Emotion Recognition

Yelin Kim and Emily Mower Provost
 Electrical Engineering and Computer Science
 University of Michigan
 Ann Arbor, USA

Email: yelinkim@umich.edu, emilykmp@umich.edu

Abstract—Human expressions are often ambiguous and unclear, resulting in disagreement or confusion among different human evaluators. In this paper, we investigate how audio-visual emotion recognition systems can leverage *prototypicality*, the level of agreement or confusion among human evaluators. We propose the use of a weighted Support Vector Machine to explicitly model the relationship between the prototypicality of training instances and evaluated emotion from the IEMOCAP corpus. We choose weights of prototypical and non-prototypical instances based on the maximal accuracy of each speaker. We then provide per-speaker analysis to understand specific speech characteristics associated with the information gain of emotion given prototypicality information. Our experimental results show that *neutrality*, one of the most challenging emotion to recognize, has the highest performance gain from prototypicality information, compared to other emotion classes: *Angry*, *Happy*, and *Sad*. We also show that the proposed method improves the overall multi-class classification accuracy significantly over traditional methods that do not leverage prototypicality.

Keywords—emotion recognition; human evaluator; prototypicality; ambiguity; neutrality

I. INTRODUCTION

Humans often perceive and evaluate the same emotion expressions in different ways. Audio-visual emotions that are labeled inconsistently by evaluators can lower the accuracy of automatic emotion recognition [1]. Prototypical emotions are, in general, less subtle displays of emotion and have high agreement rates from the evaluators. On the other hand, expressions that are hard to recognize often occur in the real world. Neutral expressions, for instance, which are either non-emotional [2], [3] or ambiguous [1], have long been considered to be one of the most challenging target emotions in affective computing [1], [4], [5]. In this research, we investigate methods to utilize *prototypicality*, the level of agreement or confusion among human evaluators, in audio-visual emotion recognition systems.

The goal of this work is to understand the advantages of including prototypicality in automatic emotion recognition systems. We use prototypicality information to weight training instances during emotion classification depending on their level of prototypicality: P_{prot} and $P_{non-prot}$ for prototypical and non-prototypical utterances, respectively (Figures 1 and 2). Our work demonstrates that among four classes of emotion that we tested (*Angry*, *Happy*, *Neutral*, and *Sad*), recognition

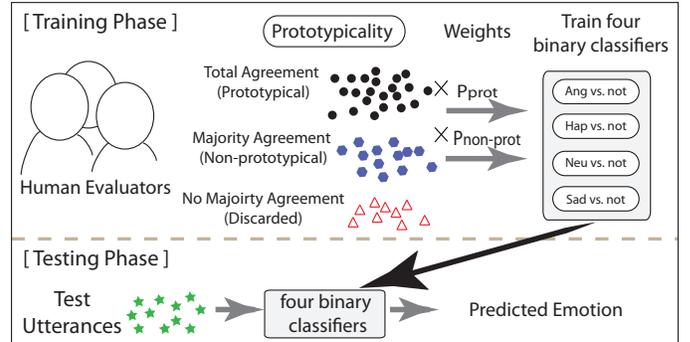


Fig. 1. Overview of our proposed system. We use *prototypicality*, the level of agreement or confusion among human evaluators, to weight training instances during emotion classification. P_{prot} and $P_{non-prot}$ represent weights for prototypical and non-prototypical utterances, respectively.

of neutral emotion achieves the highest performance gain by leveraging the prototypicality information.

There have been a few recent attempts to utilize prototypicality in emotion recognition. For instance, Eyben et al. have shown that multi-task learning of dimensional emotion labels and inter-rater standard deviation, improves the performance of dimensional emotion label regression tasks over single task learning [6]. Schuller et al. studied database and instance selection methods of eight emotion databases for training cross-corpus emotion classification systems [7]. They performed data selection by selecting prototypical training instances, where prototypicality was estimated based on an objective measure. The results showed that this method was effective for arousal recognition, but not for valence recognition. Our work differs from the previous studies, in that we provide an in-depth study about how prototypicality information is leveraged to specific emotion classes as well as certain speech characteristics.

Our proposed system first extracts audio-visual features of expressions at the utterance level. We use multi-modal feature learning, implemented using a Deep Belief Network (DBN) to generate a high-level representation of the features, shown to be effective in our previous work [8]. We use the learned features as input for binary Support Vector Machine (SVM) classifiers (e.g., angry vs. others). We weight training instances based on their prototypicality (weighted SVMs, [9]) during training, with speaker-specific and emotion-class-

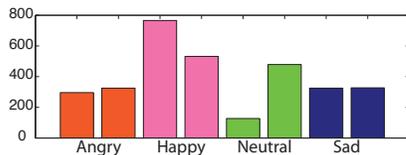


Fig. 2. The emotion class distribution of the data. The left and right bars in each emotion class represent the number of prototypical and non-prototypical utterances in that class, respectively.

specific weights. This approach provides insight into: (1) the differences in the performances of different emotion classes, focused on prototypicality and (2) speech characteristics that are correlated with system improvement due to the inclusion of prototypicality.

The novelty of our work is the demonstration of how prototypicality information can improve multi-class emotion recognition systems. We show that the use of prototypicality information leads to better system-level classification accuracy, while the accuracy improvement is the highest for neutrality recognition. We investigate how neutrality and other emotion classes are related to the prototypicality of the data. We use weights chosen based on the maximal accuracy of each speaker (oracle method) and provide per-speaker analysis to show that certain speech characteristics can be a good indicator of the potential performance gain. In our future work, we will develop automatic weighting strategies using speaker-dependent characteristics.

II. BACKGROUND

A. Related Work

1) *Prototypicality in Automatic Emotion Recognition*: Recent research in automatic emotion recognition focuses on non-prototypical and ambiguous data rather than prototypical and explicit emotions, due to its prevalence and resemblance to the real world [10]. This new paradigm has sparked interest into non-prototypical data, which contain some amount of ambiguity in human perception.

Mower et al. [1] have pointed out the importance of a system’s ability to interpret non-prototypical emotions. The authors suggested that a soft labeling scheme, instead of hard labeling, could be used to represent the degree of (estimated) presence of complex emotional expressions [11]. Lee et al. proposed a hierarchical system to recognize emotion from speech, using a decision-tree framework [12]. Their approach minimized error propagation by placing the most distinguishable emotion classes at the root of the tree, and placing more ambiguous (i.e., neutral) emotion classes in the later stage of the tree.

Eyben et al. [6] utilized prototypicality, represented as inter-rater standard deviation, by taking prototypicality as a task for multi-task learning. The authors considered five dimensional labels as targets: activation, expectation, intensity, power/dominance, and valence. Their experimental results showed that multi-task learning improved overall recognition compared to the single task modeling, achieving correlation coefficient up to 0.81 for the activation dimension and up to 0.58 for the valence dimension.

Schuller et al. [7] studied cross-corpus emotion recognition of eight widely-used emotion databases. They investigated strategies for database and instance selection using estimated prototypicality, which is calculated by the Euclidean distance from the opposite emotion class center for binary valence and arousal. The experimental results demonstrated that estimating and selecting prototypical training instances based on the distance can achieve up to 7.5% unweighted accuracy improvement in cross-corpus arousal recognition. However, the performance gain in valence recognition was not significant.

2) *Audio-Visual Emotion Recognition*: Human expressions are multi-modal, where each modality either enforces or contradicts the gestalt emotional message. Therefore, it is important for emotion recognition systems to learn the complex interactions between multiple modalities. In our previous work [8], we investigated the use of deep learning techniques in emotion recognition to learn high-level complex, non-linear interactions between audio and visual features in an unsupervised manner. We proposed a suite of deep belief networks that combined the multimodal deep learning technique proposed by Ngiam et al. [13] and feature selection methods. Our results demonstrated an improvement of audio-visual emotion recognition accuracy with deep networks. Our results showed that complex, non-linear high-level interactions between audio and visual data can be captured by deep networks to generate features that can lead to more accurate emotion recognition.

B. Audio-Visual IEMOCAP Data

In this work, we utilize the IEMOCAP database [14], an audio-visual database consisting of five female-male pairs of actors interacting given a variety of hypothetical situations and emotionally evocative scripts. We use a subset of emotions *Angry*, *Happy* (merged with the class of *Excited*), *Neutral*, *Sad*, to be consistent with previous work [11], [15], [16]. Within the dataset, there are at least three human evaluators who independently labeled each of the utterances. The inter-rater agreement provides prototypicality information: whether it has a total agreement (prototypical) or only majority agreement (non-prototypical). We remove utterances without majority vote agreement. This results in 3177 utterances in total: 621 *Angry*, 1298 *Happy*, 606 *Neutral*, and 652 *Sad* utterances. Each speaker has 62.1 ± 27.85 *Angry*, 129.8 ± 28.43 *Happy*, 60.6 ± 20.17 *Neutral*, and 65.2 ± 24.68 *Sad* utterances on average. In addition, each speaker has 151.40 ± 33.11 prototypical utterances and 166.30 ± 36.39 non-prototypical utterances on average.

The IEMOCAP database provides 3-D motion capture data that include positions of 53 markers, recorded at 120 frames per second. We use a subset of 46 markers, to be consistent with previous studies [17]. In addition to the motion capture data, we also use speech data. We calculate utterance-level static features of visual (marker positions of all facial regions) and audio (pitch, energy, Mel filter bank [2]) signal cues, similar to our previous work [8]. The eight statistical functionals used are: mean, standard deviation, quantile range, upper quantile, lower quantile, and three polynomial coefficients

TABLE I

THE FOUR-CLASS ENTROPY AND INFORMATION GAIN FOR EACH SPEAKER AND AVERAGED OVER ALL TEN SPEAKERS.

Speaker	$H(E)$	$IG(E P)$
1	1.959	0.169
2	1.905	0.187
3	1.867	0.174
4	1.853	0.061
5	1.880	0.017
6	1.960	0.047
7	1.547	0.096
8	1.910	0.036
9	1.905	0.073
10	1.774	0.091
All	1.913	0.058

with degree of three. This results in 1104 visual features and 232 audio features. We use speaker-specific z-normalization to reduce speaker-related variability in the features.

C. Preliminary Findings

We use *information gain* to understand how strongly prototypicality contributes to our ability to recognize emotion. More specifically, the information gain of a discrete random variable E (emotion) given another discrete random variable P (prototypicality) is defined as

$$IG(E|P) = H(E) - H(E|P), \quad (1)$$

where $H(E)$ denotes the entropy

$$H(E) = \sum_{e \in E} -\Pr(E = e) \log \Pr(E = e), \quad (2)$$

and $H(E|P)$ denotes the conditional entropy

$$H(E|P) = \sum_{p \in P} \Pr(P = p) H(E|P = p). \quad (3)$$

$IG(E|P)$ represents the average number of bits that one would save when transmitting the information about E given the knowledge about P . In this work, the term E describes four-emotion classes of *Angry* ('Ang'), *Happy* ('Hap'), *Neutral* ('Neu'), and *Sad*. Prototypicality P is an indicator variable, 0 if the data is prototypical (total agreement by human evaluators) and 1 if non-prototypical (majority agreement). E and P are determined from the annotated data.

Tables I and II present the entropy and information gain of emotion given knowledge of prototypicality for each speaker and averaged over all ten speakers. Table I is presented to demonstrate the overall information gain of multi-class emotion from prototypicality, whereas Table II is presented to show the specific gain for each emotion class. Table I presents the multi-class emotion recognition problem. Table II presents the binary emotion recognition problem: *Ang* vs. *Others* (everything else), *Hap* vs. *Others*, *Neu* vs. *Others*, and *Sad* vs. *Others*, respectively.

As shown in Table II, on average, we found a higher information gain for the neutral class compared to other

TABLE II

THE BINARY ENTROPY AND INFORMATION GAIN FOR EACH SPEAKER AND AVERAGED OVER ALL TEN SPEAKERS FOR THE EMOTION CLASSES OF: ANGER ('A'), HAPPINESS ('H'), NEUTRALITY ('N'), AND SADNESS ('S').

S	$H(A)$	$IG(A P)$	$H(H)$	$IG(H P)$	$H(N)$	$IG(N P)$	$H(S)$	$IG(S P)$
1	0.845	0.014	0.902	0.058	0.811	0.154	0.632	0.000
2	0.665	0.009	0.964	0.004	0.637	0.181	0.851	0.026
3	0.571	0.002	0.987	0.125	0.724	0.120	0.781	0.003
4	0.558	0.012	0.991	0.056	0.775	0.022	0.719	0.002
5	0.549	0.000	0.938	0.003	0.674	0.017	0.926	0.001
6	0.783	0.023	0.938	0.004	0.735	0.028	0.735	0.006
7	0.936	0.027	1.000	0.078	0.273	0.036	0.440	0.006
8	0.782	0.005	0.960	0.001	0.811	0.021	0.573	0.018
9	0.561	0.000	0.962	0.036	0.795	0.064	0.800	0.001
10	0.617	0.017	1.000	0.091	0.654	0.019	0.654	0.019
All	0.713	0.000	0.976	0.026	0.703	0.052	0.732	0.000

emotion classes, although the information gain varied across speakers. A comparison between Tables I and II shows that the information gain associated with the knowledge of prototypicality is similar in the four-class (Table I) and the binary neutral-class (Table II) conditions, resulting in 0.058 and 0.052 bits, respectively. We leverage the finding that different emotion classes have different information gains, by employing different weighting strategies for individual binary emotion classifiers, as described in Section III-B.

III. PROPOSED METHOD

We first use DBN to generate audio-visual features, and then build binary SVMs for emotion classification based on these features. During training of the SVMs, we assign different weights for each training instance based on its prototypicality, i.e., a *weighted SVM* strategy [9], based on prototypicality.

A. Unsupervised Feature Learning

DBN is a type of deep neural network that is built by stacking multiple restricted Boltzmann machines (RBMs). The absence of connections between units within each layer makes training of the RBM efficient [18]. For further details of DBN, see [19].

We use a two-layer DBN, similar to [8], [13], as follows: the first layer of the DBN includes two Gaussian-RBMs, each for audio and visual features, with sparsity regularization introduced in [20]. The second layer is a binary-RBM, where the concatenation of the posteriors $P(h|v)$ of each audio and video RBM is used as the observed vector for the second layer. The first Gaussian RBM uses the following energy function, given a set of observation vectors (features of training instances) $\{\mathbf{v}^{(1)}, \dots, \mathbf{v}^{(m)}\}$:

$$E(\mathbf{v}, \mathbf{h}) = \frac{1}{2\sigma^2} \sum_i v_i^2 - \frac{1}{\sigma^2} \left(\sum_i c_i v_i + \sum_j b_j h_j + \sum_{i,j} v_i W_{ij} h_j \right) + \lambda \sum_{j=1}^K \left| p - \frac{1}{m} \sum_{l=1}^m \mathbb{E} \left[h_j^{(l)} | \mathbf{v}^{(l)} \right] \right|^2, \quad (4)$$

where $\mathbb{E}[\cdot]$ is the conditional expectation given \mathbf{v} , λ is a regularization parameter, and p is a constant that specifies the target activation of the hidden unit h_j [20]. We use a sigmoid function as an activation function of each node. We fix the number of hidden nodes as 100 and 250 for each of the audio and video RBMs at the first layer, respectively, and 200 for the second layer.

B. SVM Weighting using Prototypicality of Data

For a given emotion classification task, we use a binary weighted SVM, proposed in [9] implemented in LIBSVM [21]. Each classifier uses a radial basis function kernel with a complexity parameter of 1 and a gamma of 2^{-10} , similar to [22]. Each training instance or an utterance consists of $\{x_i, t_i, P_i\}$, where x_i is the DBN-based feature vector of the i -th training instance, $t_i \in \{\text{Ang, Hap, Neu, Sad}\}$ is the emotion label of the instance, and P_i is either P_{prot} if the instance has total agreement (prototypical) or $P_{non-prot}$ if the instance has only majority agreement (non-prototypical). We solve the following modified objective function of SVM and learn w for inference:

$$\min_w \frac{1}{2} \mathbf{w}^t \mathbf{w} + C \sum_{i=1}^l P_i \xi_i. \quad (5)$$

The weight P_i of the i -th training instance in Equation 5 is assigned based on its prototypicality. An utterance is assigned a weight of P_{prot} if it is prototypical and a weight of $P_{non-prot}$ otherwise. These values are chosen based on the maximal accuracy of each speaker, shown in Table VI. We choose weights in a range of $\{1, 5, 10, 15, 20\}$ based on a speaker’s maximal accuracy.

During inference, we infer the emotion label of a j -th test instance $\{x_j\}$ using the normal vector of the SVM hyperplane \mathbf{w} , learned automatically during training. We use the trained SVMs that give the maximal accuracy for each test speaker, in an oracle manner (described in more details in Section IV-A). We use the distance from hyperplane of each emotion classifier as an estimated measure of confidence in class membership. We z-normalize the distances of each binary classifier for a balanced comparison between the outputs. Finally, we choose the emotion label with the highest distance from hyperplane, as proposed in [11]. For instance, if a test utterance has a class membership of $\{-1, 1, -1, 1\}$ (-1 for absence and 1 for presence) for each binary classifier of Ang, Hap, Neu, and Sad, with the distance from hyperplane $\{0.8, 0.75, 0.2, 0.3\}$; we then multiply each class membership with the corresponding distances, obtaining $\{-0.8, 0.75, -0.2, 0.3\}$. The final emotion prediction of the utterance is the class of happiness.

We use leave-one-speaker-out cross validation, where we use one held-out speaker for testing, and the remaining nine speakers for training. To be consistent with previous work on the IEMOCAP dataset, we report an average of the four-class recall.

IV. RESULTS AND DISCUSSION

In this section, we present four-class emotion recognition experiments using prototypicality information. We report

TABLE III
BASELINE (NO WEIGHTING) CLASSIFICATION ACCURACY FOR EACH SPEAKER AND OVER ALL SPEAKERS. THE COLUMNS REPRESENT PER-CLASS PERFORMANCE AND THE UNWEIGHTED ACCURACY (UW) OVER ALL FOUR CLASSES ($A = \text{angry}$, $H = \text{happy}$, $N = \text{neutral}$, $S = \text{sad}$).

Speaker	UW	A	H	N	S
1	62.76	61.11	79.76	36.36	73.81
2	63.14	87.72	73.44	32.08	59.34
3	62.17	74.29	69.64	48.08	56.67
4	72.33	85.00	60.29	57.14	86.89
5	64.99	84.21	77.36	41.51	56.86
6	69.60	82.65	70.47	45.98	79.31
7	54.97	52.50	64.94	25.00	77.42
8	60.47	60.00	60.75	34.29	86.84
9	73.78	78.26	84.44	55.95	76.47
10	70.59	80.00	56.89	60.00	85.45
Mean	65.48	74.57	69.80	43.64	73.91
Std	5.95	12.31	9.11	11.80	12.10

TABLE IV
CLASSIFICATION ACCURACY WHEN WEIGHTING INDIVIDUAL BINARY CLASSIFIERS (E.G., W_A IS THE SYSTEM WHEN USING A WEIGHTED ANGRY SVM). THE COLUMNS REPRESENT UNWEIGHTED ACCURACY OF THE SYSTEM FOR EACH SPEAKER AND OVER ALL TEN SPEAKERS ($A = \text{angry}$, $H = \text{happy}$, $N = \text{neutral}$, $S = \text{sad}$).

Speaker	W_A	W_H	W_N	W_S
1	63.72	63.79	64.43	62.76
2	63.81	64.98	66.53	63.22
3	63.82	62.84	65.55	64.12
4	72.33	74.24	73.40	72.33
5	66.35	66.11	64.99	67.42
6	70.55	69.77	70.18	70.04
7	55.03	55.83	57.03	55.41
8	60.47	61.09	60.47	61.01
9	74.52	73.97	75.29	74.86
10	71.19	72.29	72.40	70.84
Mean	66.18	66.49	67.03	66.20
Std	6.02	6.02	5.81	5.95

unweighted and per-class accuracy. We use a paired t-test method for leave-one-speaker-out cross validation to test the significance level between methods, as suggested in [23]. The baseline method of our experiments is a traditional SVM that does not utilize any prototypicality information in emotion recognition (Table III). In the first experiment (Section IV-A, Tables IV and V), we use weighting strategies for one of the four binary emotion classifiers, while the other three classifiers use a traditional method without any weighting strategies. This is to explore the utility of prototypicality in recognition of individual emotions. Our hypothesis is that an emotion class which has higher information gain, neutrality for instance, will result in a higher performance gain as well, when using the weighting strategy. In the second experiment (Section IV-A, Table VI), we combine the weights of each binary classifier resulting in the per-speaker maximal accuracy (Table VII) and

TABLE V

CLASSIFICATION ACCURACY WHEN WEIGHTING INDIVIDUAL BINARY CLASSIFIERS (E.G., W_A IS THE SYSTEM WHEN USING A WEIGHTED ANGRY SVM). THE COLUMNS REPRESENT THE ACCURACY OF EACH OF THE FOUR EMOTION CLASSES AS A FUNCTION OF WEIGHTING THE INDIVIDUAL BINARY SVMs. FOR EXAMPLE, THE COLUMN UNDER THE HEADERS ‘‘ANGRY’’ AND ‘‘ W_A ’’ IS THE ACCURACY OF THE ANGRY PREDICTION GIVEN A WEIGHTED ANGRY BINARY SVM.

S	Angry				Happy				Neutral				Sad			
	W_A	W_H	W_N	W_S	W_A	W_H	W_N	W_S	W_A	W_H	W_N	W_S	W_A	W_H	W_N	W_S
1	65.28	62.50	63.89	59.72	80.95	80.95	78.57	79.76	34.85	37.88	37.88	34.85	73.81	73.81	73.81	73.81
2	87.72	87.72	87.72	87.72	74.22	78.91	72.66	73.44	32.08	33.96	39.62	32.08	59.34	59.34	64.84	61.54
3	80.00	74.29	74.29	74.29	69.64	72.32	66.96	68.75	46.15	48.08	50.00	46.15	56.67	56.67	55.00	61.67
4	85.00	85.00	85.00	85.00	59.56	65.44	58.82	58.82	55.71	57.14	61.43	57.14	86.89	86.89	86.89	86.89
5	86.84	86.84	84.21	84.21	76.42	79.25	71.70	79.25	41.51	41.51	41.51	41.51	56.86	56.86	54.90	64.71
6	83.67	82.65	82.65	81.63	70.47	71.14	69.13	69.13	49.43	45.98	47.13	45.98	79.31	79.31	80.46	82.76
7	53.33	53.33	51.67	52.50	64.94	68.39	64.37	63.22	25.00	25.00	31.25	25.00	77.42	77.42	77.42	80.65
8	60.00	60.00	60.00	60.00	60.75	61.68	58.88	59.81	34.29	34.29	34.29	34.29	86.84	86.84	84.21	84.21
9	78.26	78.26	82.61	78.26	83.70	85.19	85.93	84.44	55.95	55.95	57.14	57.14	69.41	76.47	76.47	77.65
10	80.00	80.00	78.00	78.00	55.69	58.08	56.89	56.89	54.55	61.82	67.27	58.18	85.45	85.45	85.45	89.09
Mean	76.01	75.06	75.00	74.13	69.63	72.14	68.39	69.35	42.95	44.16	46.75	43.23	73.20	73.91	73.95	76.30
Std	12.10	12.24	12.32	12.33	9.37	8.89	9.27	9.70	11.06	11.80	12.05	11.72	12.14	12.10	11.88	10.40

use the weighting strategies for all of the emotions.

Finally, we analyze the relationship between the information gain associated with prototypicality and speaker-dependent speech features (Section IV-A, Table VIII). This analysis will provide insight into future work that can automatically decide weighting strategies for individual speakers based on the specific audio-visual characteristics.

A. Weighting Scheme for Individual Emotions

Tables IV and V show the unweighted four-class and per-class accuracy results of each speaker, respectively, when weighting schemes are used for only one of the four binary emotion classifiers: *Ang vs. Others*, *Hap vs. Others*, *Neu vs. Others*, and *Sad vs. Others*. This results in four weighting strategies, depending on which binary emotion classifier employs the weighted SVM: Ang (W_A scheme), Hap (W_H scheme), Neu (W_N scheme), and Sad (W_S scheme). The maximal accuracy was chosen for each of the schemes, for individual speakers (chosen weights shown in Table VII).

The improvement in four-class accuracy from using the weighting scheme (Table IV) compared to the baseline (Table III) is greatest for W_N , i.e., we only weight the data for *Neu vs. Others* binary classifier and use traditional SVMs for the other three binary classifiers, with an increase by 1.55% (65.48% to 67.03%, $p < 0.005$, significant). The weighting for the other schemes results in 0.70% improvement for W_A (65.48% to 66.18%, $p < 0.005$, significant), 1.01% for W_H (65.48% to 66.49%, $p < 0.001$, significant), and 0.72% for W_S (65.48% to 66.20%, $p < 0.02$, significant).

A comparison between per-class accuracy results in Tables III (baseline) and V (weighting) shows the significant improvement of all four emotion classes by using the weighting strategy compared to the baseline. The per-class angry accuracy of W_A compared to the baseline improves by 1.44% (74.57% to 76.01%, $p < 0.05$, significant). The improvements using the other schemes are 2.34% for the per-class happy accuracy using W_H (69.80% to 72.14%, $p < 0.02$, significant), 3.11% for the per-class neutral accuracy (43.64% to 46.75%, $p < 0.01$, significant), 2.71% for the per-class sad accuracy (73.91%

TABLE VI

CLASSIFICATION ACCURACY WHEN USING A SYSTEM COMPOSED OF OPTIMALLY WEIGHTED SVMs FOR EACH OF THE FOUR EMOTION CLASSES. THE COLUMNS REPRESENT PER-CLASS PERFORMANCE AND THE UNWEIGHTED ACCURACY (UW) OVER ALL FOUR CLASSES ($A = angry$, $H = happy$, $N = neutral$, $S = sad$).

Speaker	UW	A	H	N	S
1	64.78	65.28	79.76	37.88	76.19
2	66.05	87.72	71.09	33.96	71.43
3	65.30	88.57	64.29	50.00	58.33
4	74.95	87.50	63.97	61.43	86.89
5	68.31	86.84	76.42	45.28	64.71
6	69.37	81.63	67.11	48.28	80.46
7	59.52	55.00	61.49	31.25	90.32
8	61.62	61.54	60.75	40.00	84.21
9	74.22	78.26	81.48	57.14	80.00
10	71.79	80.00	58.08	63.64	85.45
Mean	67.59	77.23	68.44	46.89	77.80
Std	5.12	12.24	8.32	11.31	10.24

to 76.30%, $p < 0.03$, significant). This demonstrates that the emotion recognition improves significantly by employing prototypicality, where the greatest improvement is achieved for neutrality recognition.

B. Combined Weighting Scheme

We combine the weighting schemes that yielded the maximal accuracy for each emotion class in Section IV-A. As shown in Table VI, the maximal accuracy achieved by combining these weights for each emotion is 67.59%, an improvement of 2.11% compared to the system that does not use prototypicality weighting ($p < 0.001$, significant).

The results suggest that the use of prototypicality information helps the overall multi-class recognition. We also found that the multi-class performance gain using only W_N (67.03%) can achieve comparable performance ($p = 0.2$) compared to a system that utilizes the weighting strategy for all four emotion classifiers (67.59%). This suggests that the recognition of

TABLE VII
OPTIMAL WEIGHTS FOR EACH OF THE EMOTION CLASSES.

Speaker	Angry		Happy		Neutral		Sad	
	prot	non-prot	prot	non-prot	prot	non-prot	prot	non-prot
1	20	5	10	5	20	1	1	1
2	5	1	5	1	20	10	5	1
3	5	15	20	1	15	1	1	5
4	1	1	20	10	5	5	1	1
5	20	1	20	1	1	1	20	1
6	5	1	15	10	5	1	5	1
7	15	1	5	1	1	20	20	1
8	1	1	10	5	1	1	1	5
9	1	20	5	1	10	1	15	1
10	20	20	20	5	20	5	10	1

neutrality benefits most from employing weighting based on prototypicality. Compared to the state-of-the-art results on the IEMOCAP dataset (66.12%) that is achieved in our previous work [8], the highest accuracy we obtained in our paper is higher (67.59%, relative improvement of 1.47%).

C. Analysis on Speaker-dependent Speech Characteristics

Table VIII shows that speakers 1, 2, and 3 have the highest information gain and that speakers 5, 8, and 10 have the lowest information gain. We analyze the speech characteristics of each speaker to understand why certain speakers have higher benefits from prototypicality compared to others (Table VIII). We present a subset of pitch and energy features that are correlated (Pearson product-moment correlation, values greater than 0.5 or less than -0.5) with speaker’s information gain of ground truth neutrality given the prototypicality information, $IG(Neu|Prot)$, in Table VIII: the bottom row of the table shows the correlation between $IG(Neu|Prot)$ and (i) mean pitch of neutral, (ii) mean energy of neutral, (iii) mean energy of happy, (iv) mean energy of non-prototypical, utterances, as well as the (iv) standard deviation energy of neutral and (v) standard deviation energy of non-prototypical utterances. The results show strong negative correlations between information gain and mean pitch, as well as between information gain and mean energy of neutral utterances, both of -0.66. In other words, the mean pitch and energy of neutral data tend to be smaller for speakers with higher information gain. Table VIII also shows that the standard deviation of energy features for neutral utterances tend to be smaller (correlation -0.74). On the other hand, the standard deviation of energy features for non-prototypical (of all four emotion classes) utterances are higher (correlation 0.73) for the speakers with higher information gain. This may indicate that speakers, with consistently low mean and low standard deviation in their energy features during their neutral states, achieve higher benefit from prototypicality for neutrality recognition.

We also highlight the strong correlation of 0.74 between per-speaker performance gain of the neutral classifier (UW accuracy of 67.03% for W_N scheme in Table IV) and per-

TABLE VIII
THE CORRELATION BETWEEN SPEAKER CHARACTERISTICS AND $IG(Neu|Prot)$. THE RESULTS ARE PRESENTED OVER EACH SPEAKER AND AVERAGED OVER ALL SPEAKERS. THE COLUMN ‘NEU’ DENOTES NEUTRAL UTTERANCES, ‘HAP’ DENOTES HAPPY UTTERANCES, AND ‘NON-PROT’ DENOTES NON-PROTOTYPICAL UTTERANCES.

S	$IG(Neu Prot)$	mean pitch	mean energy			std energy	
		Neu	Neu	Hap	Non-prot	Neu	Non-prot
1	0.154	-0.448	-0.535	0.230	-0.277	0.702	0.874
2	0.181	-0.193	-0.241	0.177	-0.103	0.839	0.924
3	0.120	-0.140	-0.055	0.165	-0.192	0.829	0.971
4	0.022	-0.096	0.101	0.219	-0.046	0.856	0.984
5	0.017	-0.136	0.058	0.178	-0.162	0.938	0.893
6	0.028	-0.060	0.076	-0.003	-0.088	0.921	0.916
7	0.036	-0.197	-0.321	0.032	-0.155	0.878	0.995
8	0.021	-0.069	-0.055	0.110	-0.008	0.907	0.974
9	0.064	-0.157	-0.031	0.158	-0.097	0.892	0.941
10	0.019	-0.143	0.004	0.035	-0.121	0.878	1.001
Corr		-0.66	-0.66	0.51	-0.50	-0.74	0.73

speaker information gain in neutral emotions given prototypicality ($IG(Neu|Prot)$ column in Table II). For other emotion classes, the correlation was not significant (absolute value less than 0.5), which indicates stronger correlation between neutrality recognition and $IG(Neu|Prot)$. These comparisons demonstrate that the use of prototypicality labels that are highly associated with emotion classes results in higher performance gain.

V. CONCLUSIONS AND FUTURE WORK

In this paper, we explore the use of prototypicality information (i.e., inter-rater agreement level) in audio-visual emotion recognition. We found that the neutral emotion class has the highest information gain as well as performance gain in multi-class emotion recognition tasks by utilizing prototypicality information. We employed a weighted SVM method proposed in [9] to validate our hypothesis, and demonstrated that the best accuracy at 67.59% (1.47% higher than the state-of-the-art result [8]) is achieved when prototypicality information is used. We found that this is significantly higher than traditional baseline methods that do not employ any prototypicality information in the system. The results are also comparable when the weighting is used only on neutrality recognition and when the weighting is used in all four emotion recognition (paired t-test, $p = 0.2$). This suggests that the weighting strategy improves overall recognition accuracy by improving the neutral recognition task.

The limitation of our training method is that the system requires weights chosen based on the maximal accuracy of each speaker. In future work, we will develop automatic weighting strategies based on speaker-dependent characteristics. Also, a remaining open question is the working definition of prototypicality in emotion recognition. Future work will explore different representation strategies of non-prototypicality, e.g., ‘soft’ mixtures of multiple prototypes [11].

REFERENCES

- [1] E. Mower, A. Metallinou, C.-C. Lee, A. Kazemzadeh, C. Busso, S. Lee, and S. Narayanan, "Interpreting ambiguous emotional expressions," in *3rd International Conference on Affective Computing and Intelligent Interaction and Workshops*.
- [2] C. Busso, S. Lee, and S. Narayanan, "Using neutral speech models for emotional speech analysis," *Proceedings of Interspeech*, pp. 2225–2228, 2007.
- [3] A. Batliner, S. Steidl, and E. Nöth, "Releasing a thoroughly annotated and processed spontaneous emotional database: the FAU Aibo emotion corpus," in *Proc. of a Satellite Workshop of LREC*, 2008, pp. 28–31.
- [4] F. Burkhardt, M. Van Ballegooy, K.-P. Engelbrecht, T. Polzehl, and J. Stegmann, "Emotion detection in dialog systems: applications, strategies and challenges," in *3rd International Conference on Affective Computing and Intelligent Interaction and Workshops*.
- [5] Y. Kim and E. Mower Provost, "Emotion classification via utterance-level dynamics: A pattern-based approach to characterizing affective expressions," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Vancouver, BC, May 2013, pp. 3677–3681.
- [6] F. Eyben, M. Wöllmer, and B. Schuller, "A multitask approach to continuous five-dimensional affect sensing in natural speech," *ACM Transactions on Interactive Intelligent Systems (TiiS)*, vol. 2, no. 1, p. 6, 2012.
- [7] B. Schuller, Z. Zhang, F. Wenginger, and G. Rigoll, "Selecting training data for cross-corpus speech emotion recognition: Prototypicality vs. generalization," in *Proc. 2011 Afeka-AVIOs Speech Processing Conference, Tel Aviv, Israel*. Citeseer, 2011.
- [8] Y. Kim, H. Lee, and E. Mower Provost, "Deep learning for robust feature generation in audiovisual emotion recognition," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Vancouver, BC, March 2013, pp. 3687–3691.
- [9] X. Yang, Q. Song, and Y. Wang, "A weighted support vector machine for data classification," *International Journal of Pattern Recognition and Artificial Intelligence*, vol. 21, no. 05, pp. 961–976, 2007.
- [10] B. Schuller, S. Steidl, A. Batliner, F. Burkhardt, L. Devillers, C. Müller, and S. Narayanan, "Paralinguistics in speech and language—state-of-the-art and the challenge," *Computer Speech & Language*, vol. 27, no. 1, pp. 4–39, 2013.
- [11] E. Mower, M. Mataric, and S. Narayanan, "A framework for automatic human emotion classification using emotion profiles," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 5, pp. 1057–1070, 2011.
- [12] C.-C. Lee, E. Mower, C. Busso, S. Lee, and S. Narayanan, "Emotion recognition using a hierarchical binary decision tree approach," in *Interspeech*, Jeju Island, South Korea, October 2009, pp. 320–323.
- [13] J. Ngiam, A. Khosla, M. Kim, J. Nam, H. Lee, and A. Ng, "Multi-modal deep learning," in *International Conference on Machine Learning (ICML)*, 2011, pp. 689–696.
- [14] C. Busso, M. Bulut, C. Lee, A. Kazemzadeh, E. Mower, S. Kim, J. Chang, S. Lee, and S. Narayanan, "Temocap: Interactive emotional dyadic motion capture database," *Language Resources and Evaluation*, vol. 42, no. 4, pp. 335–359, 2008.
- [15] A. Metallinou, C. Busso, S. Lee, and S. Narayanan, "Visual emotion recognition using compact facial representations and viseme information," in *IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP)*, Dallas, TX, March 2010, pp. 2474–2477.
- [16] S. Mariosyad and C. Busso, "Feature and model level compensation of lexical content for facial emotion recognition," in *IEEE International Conference on Automatic Face and Gesture Recognition (FG)*, Ljubljana, Slovenia, May 2013.
- [17] Y. Kim and E. Mower Provost, "Say cheese vs. smile: Reducing speech-related variability for facial emotion recognition," in *Proceedings of the ACM International Conference on Multimedia (ACM MM'14)*, 2014.
- [18] G. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [19] Y. Bengio, "Learning deep architectures for AI," *Foundations and Trends in Machine Learning*, vol. 2, no. 1, pp. 1–127, 2009.
- [20] H. Lee, C. Ekanadham, and A. Ng, "Sparse deep belief net model for visual area v2," *Advances in Neural Information Processing Systems (NIPS)*, vol. 20, pp. 873–880, 2008.
- [21] C. Chang and C. Lin, "Libsvm: a library for support vector machines," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 2, no. 3, p. 27, 2011.
- [22] Y. Shanguan and E. Mower Provost, "Emoshapelets: Capturing local dynamics of audiovisual affective speech," in *Affective Computing and Intelligent Interaction (ACII)*, Xi'an, China, September 2015.
- [23] T. G. Dietterich, "Approximate statistical tests for comparing supervised classification learning algorithms," *Neural Computation*, vol. 10, no. 7, pp. 1895–1923, 1998.