

Influence of Different Feature Selection Approaches on the Performance of Emotion Recognition Methods Based on SVM

Daniil Belkov, Konstantin Purtov, Vladimir Kublanov
 Ural Federal University (UrFU)
 Yekaterinburg, Russia
 d.d.belkov, k.s.purtov@gmail.com, kublanov@mail.ru

Abstract—In this paper we evaluate performance of modern emotion recognition methods. Our task is to classify emotions as basic 8 categories: anger, contempt, disgust, fear, happy, sadness, surprise and neutral. CK+ dataset is used in all experiments. We apply Adaptive Boosting and Principal Component Analysis for dimensionality reduction and Support Vector Machine for classification. Size of train dataset is increased by use of few frames of sequences instead of one and vertical mirroring of faces. All images were normalized with mean centering and standardizing. In total 4428 images were used in experiment. The proposed method can work in real time and achieved average accuracy higher than 95%.

I. INTRODUCTION

Human's facial expressions provide variable information about emotions and internal state of the person. Ability to automatically recognize facial expressions offers vast possibilities for video surveillance systems, systems that measure audience mood, public security and control systems, etc. Vision-based capture systems attempt to provide such applications, by using video cameras as remote sensors.

Video based emotion analysis is a very challenging problem. In past decade, a large number of systems for recognizing emotion [1], [2], [3] were developed with the use of modern methods of machine learning and high-quality databases, such as CK + [4], MMI [5], Jaffe [6].

These systems interpret the expression as one of the seven basic emotions (happiness, anger, contempt, disgust, sadness, surprise and fear). The interpretation is based on the analysis of facial expression using Facial Action Coding System (FACS) [7].

The purpose of this article is to create a real-time system of monitoring of emotional state. To accomplish it we evaluate the accuracy of emotion classification algorithms with using the different feature sets. Feature selection was carried out by using the Principal Component Analysis (PCA) and Adaptive Boosting (AdaBoost) methods. Support Vector Machine (SVM) with Radial Basis Function (RBF) kernel was used as the learning algorithm of emotion classification.

The paper is organized as follows. Section II contains a short overview of existing applications and presents the accuracy of emotion classification. Section III gives a detailed description of used approach, including the feature selection, images preprocessing and classification methods. Section IV

describes the emotion dataset that was used in this work. Section V presents the experiment overview and evaluation metrics and next, Section VI combines all calculated results, with overview of the accuracy. All conclusions of paper presented in Section VII.

II. RELATED WORK

A. Computer Emotion Recognition Toolbox (CERT)

CERT [2] is a fully automatic, real-time software tool that estimates facial expression both in terms of 19 FACS Action Units, as well as the 6 universal emotions.

This system uses its own face detector, which was trained using an extension of the Viola-Jones [9] approach. It employs GentleBoost as the boosting algorithm and WaldBoost for automatic cascade threshold selection. On the CMU+MIT dataset, CERT's face detector achieves a hit rate of 80.6

After finding the face region, CERT estimates positions of 10 feature points: the corners of the eyes, eye centers, tip of the nose, the corners of the mouth and mouth center. Each facial feature detector, trained using GentleBoost, outputs the log-likelihood ratio of that feature being present at a location (x, y) within the face, to being not present at that location.

Given the set of 10 facial feature positions, the face patch is re-estimated at a canonical size of 96×96 pixels using an affine warp. The warp parameters are computed to minimize the L2 norm between the warped facial feature positions of the input face and a set of canonical feature point positions.

The cropped 96×96 -pixel face patch is then convolved with a filter bank of 72 complex-valued Gabor filters of 8 orientations and 9 spatial frequencies. The magnitudes of the complex filter outputs are concatenated into a single feature vector.

This feature vector is input to a separate linear support vector machine (SVM). SVM classifier provides distance of the input feature vector to the SVM's separating hyperplane, which can be interpreted as the intensity of certain facial movement.

To determine the 6 basic emotions and the neutral facial expression the classifier based on multivariate logistic regression was used. AU intensities and emotion ground truth labels were used for training classifier.

For this system the average accuracy is 0.93, the average F-measure is 0.79.

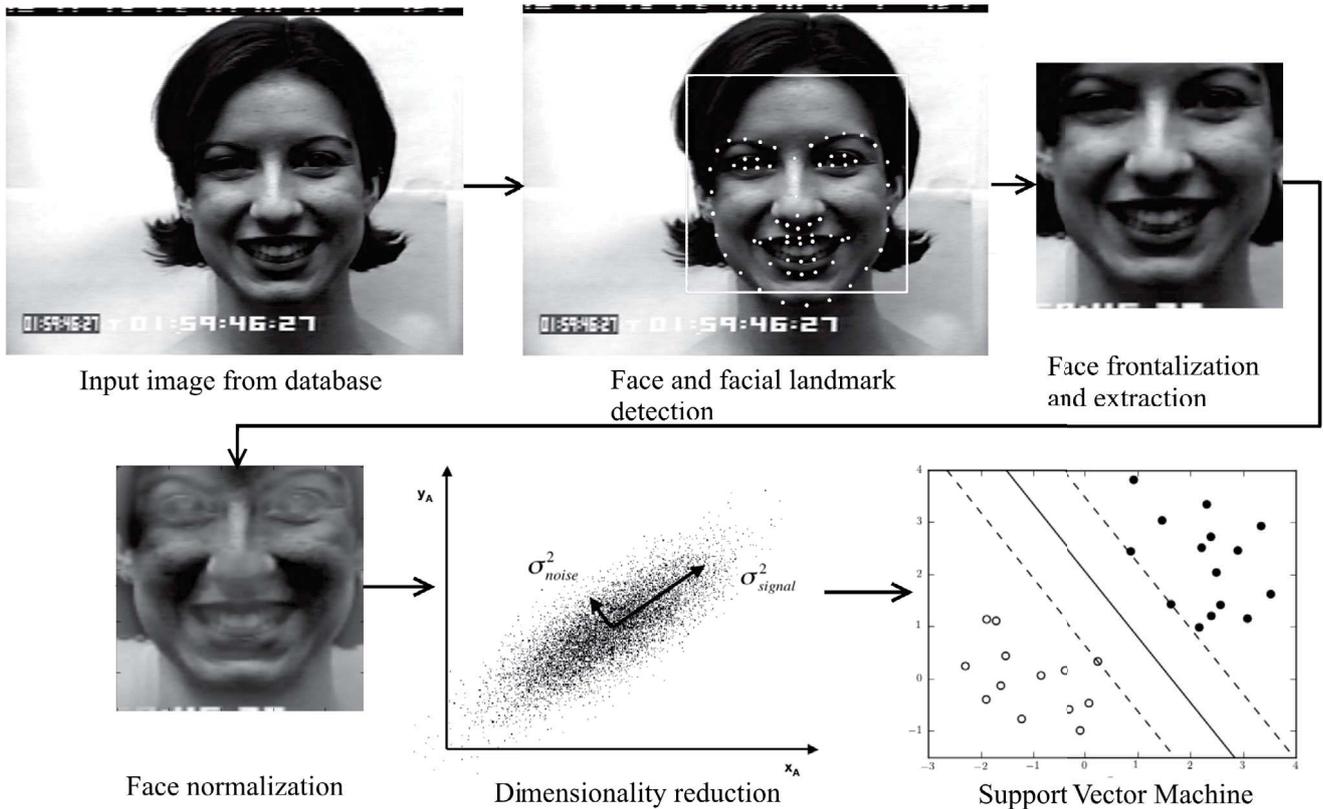


Fig. 1. Scheme of emotion recognition algorithm

B. Facial expression recognition using radial encoding of local Gabor features and classifier synthesis

In this approach, [3] mouth and eye areas are manually labeled, then the face area of size 184×152 pixels is determined by these areas. Then, each image was divided into several local regions with a 50% overlap.

Next, to each of these local regions they apply the bank of Gabor filters with 3 scales and 8 orientations. The outputs of the filters were converted into a feature matrix.

Then, each filtered image was encoded by using a radial grid. The radial grid of resolution 18×5 , with the center at the center of local region was used.

To reduce dimensionality of the feature matrix PCA and Fisher's linear discriminant were used.

For the classification of emotions K-nearest neighbor algorithm with $k = 1$ was used. The best result on the CK dataset was 91.51%.

In contrast to these methods, we do not use AU, because of the small number of labeled samples for training the robust classifier. Therefore, we take the last 20% of the frames of each sequence as samples emotions.

III. ALGORITHM

The proposed system for the classification of emotions consists of 6 main stages as shown in Fig. 1: (A) face detection,

(B) detection of the key points, (C) face frontalization and extraction, (D) normalization of face images, (E) dimensionality reduction (F) classification of emotions. Next we briefly describe each of the stages.

A. Face detection

First we have to find a face on each image in the database. We use Viola Jones approach [9]. This algorithm is widely used for face detection because it has high accuracy and can process images in real time.

B. Detection of the key points

The next stage is to find the key points in a given face region. We use face alignment tool based on the algorithm that uses cascades of regression trees [8]. This tool outputs location estimates of 68 key points: contour of the face, the nose, the outer and inner sides of the lips, eyes, eyebrows. Given the initial constellation of the (x, y) locations of the 68 facial features, the location estimates are refined using linear regression. Outputs of the key points detector marked by white circles within the face in Fig. 1.

C. Face frontalization and extraction

Given the face bounding rectangle and the key points we frontalize the face image with the affine transformation. Parameters of affine transformation are calculated to minimize

the L2 norm between the warped facial feature positions of the input face and a set of canonical feature point positions.

Then we extract face patch size of 96×96 pixels from the original image. Each patch is then presented in the form of vector length 9216.

D. Normalization of the face image

Each resulting face image is normalized using the following strategy. First, from each pixel of the image we subtract the mean value of pixels in the image, as shown in equation 1.

$$M_i = I_i - \bar{I}_i \quad (1)$$

where I_i is input image, \bar{I}_i – mean value in current image.

Then, from each pixel we subtract the average value of that pixel on all the images in the database and divide the resulting value by the standard deviation of that pixel on all the images, as shown in equation 2 :

$$X_{x,y} = \frac{M_{x,y} - \bar{M}_{x,y}}{std(M_{x,y})} \quad (2)$$

where x and y is coordinates of pixel, std – standart deviation of pixel.

E. Dimensionality reduction

1) *AdaBoost*: AdaBoost algorithm trains a cascade of weak classifiers in iterative manner modifying weights of training samples. The weights are changed according to the correctness of object classification at the current stage. If the object was classified incorrectly, then its weight increases, otherwise reduces. The next predictor in the cascade will focus more on those objects that were incorrectly classified at the previous stage.

Decision stumps or decision trees are usually used as the weak classifiers. In this paper we use decision trees with tree depth equal to 4. The number of weak classifiers in a cascade is set to 50.

Since we use the decision trees as the weak classifiers, we can evaluate the importance of each feature, i.e. in this case the importance of values of each feature vector component. The basic idea is to evaluate the significance of each feature for separating in the tree nodes. In the case of decision trees with depth greater than 1 we determine the importance of each feature to a tree, then determine the final importance by averaging these values. Features in the top of the tree have a greater weight than the features used in the bottom, as they separate more objects.

Given the obtained importance of the features we can reduce the dimensionality of space by taking into account the most informative features.

In this study, we use all the features that have the importance greater than 0. The total number of such features is 348.

2) *PCA*: PCA is a method that is used in data mining to reduce the dimensionality. Images often contain redundant information that has little importance in the analysis.

This method reduces the dimensionality of data and hence reduce the computational load during their processing, while losing a minimal amount of information.

First, training set is centered on the mean value. Then the covariance matrix is calculated.

From covariance matrix the Eigen values and eigenvectors are calculated. The eigenvector with the highest Eigen values is the principal component. Eigen values are ordered in ascending manner to form feature matrix. Eigenvectors with low Eigen values can be dropped. The principal component data set is done by multiplying transposed data set value and transposed feature vectors.

In this paper we used the PCA with number of principal components set to 348 for comparison with AdaBoost and with number of principal components set to 1800 as it provides the best results.

F. Classification of emotions

For the classification of emotions, we use a support vector machine algorithm [10]. This algorithm allows to separate data in a high dimensional space using hyperplane located at a maximum margin of all classes.

We use SVM which implements the “one-against-one” classification strategy. In this approach, the classifier decides for each pair of classes. The final classification is made by voting. The RBF function is selected as a kernel function of the classifier.

Kernel function selection is based on fact, that emotions are not linearly separable, since the same AU can be present in different emotions. For example, AU1 “inner brow raiser” presented in emotions of fear, sadness and surprise, and AU15 “lip corner depressor” presented in emotions of sadness and disgust. RBF kernel allow to separate that kind of data. It projects data to a higher dimensional feature space, where it can be separate. So RBF function is the best default choice when data is not linearly separable as it generally gives the better result than other kernels.

IV. DATABASE

CK+ dataset consists of a frame sequences. In each sequence there is a person performing a certain facial expression. Resolution of all frames is 640×480 pixels. Examples of images from CK+ dataset presented in Fig. 2. Left column represent the Neutral emotion state. Other colums contains variaes emotions for the same person.

Participants were 18 to 50 years of age, 69% female, 81%, Euro-American, 13% Afro-American, and 6% other groups. CK + contains the sequences for 123 subjects. The sequence length varies from 10 to 60 frames. Each of the sequences contains images from onset (neutral frame) to peak expression (last frame). The peak frame was reliably FACS coded for facial action units.



Fig. 2. Example of images from CK+ dataset

Emotion label is based on the subject's impression of each of the 7 basic emotion categories: anger, contempt, disgust, fear, happy, sadness and surprise.

Labels are defined according to the FACS codes and the Emotion Prediction Table from the FACS manual. In total 327 sequences have emotion labels. We use these sequences in the experiment.

V. EXPERIMENT

A. Dataset preprocessing

We divide the CK+ dataset as follows: first 6% of sequence of frames used as a neutral facial expression, last 20% used as a facial expression that represents a certain emotion. This approach allows us to increase size of dataset. After extraction all samples were reviewed manually to exclude any class overlapping.

Furthermore, each image was mirror reflected by vertical axis. Such approach provides increasing the size of dataset and increasing the robustness of the algorithm. In the end, there are 4428 face images in the dataset.

To obtain different sets we randomly shuffled and divide our image set to training and testing subsets 10 times. Subsets were divided at a ratio of 70:30. In our case, each training subset contains 3104 images and testing subset contains 1324 images. So, in test subset there are about 500 images with Neutral facial expression, and 110 for each of emotion.

B. Evaluation metrics

To evaluate performance of proposed algorithms we use Precision, Recall, Accuracy and F_1 -score. All classifier's decisions may be divided into four groups:

- TP — true positive decision,

- TN — true negative decision,
- FP — false positive decision,
- FN — false negative decision.

Accuracy can be found as:

$$Accuracy = \frac{TP + TN}{TP + TN + FP + FN} \quad (3)$$

F-measure can be found as:

$$F_1 = 2 \cdot \frac{precision \times recall}{precision + recall} \quad (4)$$

Here precision and recall are:

$$precision = \frac{TP}{TP + FP} \quad (5)$$

$$recall = \frac{TP}{TP + FN} \quad (6)$$

C. Hardware and Software

All computations were held in a personal computer with following configuration: Intel Core i7 4770, 3,4 GHz, and 8 Gb DDR4 RAM. The software was implemented in the Python 2.7, with using the popular open-source packages OpenCV 3.1, scipy, numpy, scikit-learn.

VI. RESULTS

In this section we present the results of this study. All experiments were conducted on the same dataset. Dataset was randomly shuffled and divided to train and test subsets 10 times, to obtain reliable accuracy assessment. Since the experiment was performed on 10 different subsets, we present the average results among all datasets for all algorithms used in experiment.

Results are shown in Tables I, II, III for various methods of reducing feature dimension. The values of Accuracy, F_1 -score, Precision, Recall were calculated according to equations 3,4,5,6, respectively.

Table I shows that average Accuracy for AdaBoost dimension reduction are higher than 95% for all Emotion states. The highest values of Accuracy correspond to Anger, Disgust, Fear, Happy, Surprise emotion states. The Sad emotion has a little bit worse result of accuracy. The worst values obtained for a Contempt and Neutral emotion states, 96.95% and 95.18% respectively. F_1 -score results are similar to Accuracy. For all emotion states, except Neutral, Precision value is close to 1 and higher than Recall values. It means that current emotion state was classified correctly in almost all cases, but sometimes it classified another emotions as current.

Table II shows the results for PCA dimension reduction method with first 348 principal components. In this case, Accuracy of all emotion states, except Neutral, is close to 90%. The Neutral emotions state has worst result equal to 58.88%.

TABLE I. RESULTS OF SVM CLASSIFICATION USING ADABOOST WITH 348 COMPONENTS

Emotion	Anger	Contempt	Disgust	Fear	Happy	Sad	Surprise	Neutral	Average
Anger	98.42	0	0	0	0	1.97	0	0.04	
Contempt	0	76.88	0	0	0	0	0.73	0.04	
Disgust	0.08	0	98.14	0	0	0	0	0	
Fear	0.08	0	0	100	0.23	0	0	0.04	
Happy	0	0	0	0	99.77	0	0	0	
Sad	0.15	0	0	0	0	86.76	0	0.15	
Surprise	0	0	0	0	0	0	98.65	0	
Neutral	1.28	23.13	1.86	0	0	11.27	0.62	99.74	
Precision	0.98	0.99	1.00	1.00	1.00	1.00	1.00	0.72	0.96
Recall	0.98	0.77	0.98	1.00	1.00	0.87	0.99	1.00	0.95
F1-score	0.98	0.87	0.99	1.00	1.00	0.93	0.99	0.84	0.95
Accuracy	99.53	96.95	99.74	99.96	99.97	98.25	99.82	95.18	98.67

TABLE II. RESULTS OF SVM CLASSIFICATION WITH USING FIRST 348 PCA COMPONENTS

Emotion	Anger	Contempt	Disgust	Fear	Happy	Sad	Surprise	Neutral	Average
Anger	60.60	0	0	0	0	0	0	0	
Contempt	0	50.63	0	0	0	0	0	0	
Disgust	0	0	35.85	0	0	0	0	0	
Fear	0	0	0	59.44	0	0	0	0	
Happy	0	0	0	0	69.94	0	0	0	
Sad	0	0	0	0	0	45.35	0	0.09	
Surprise	0	0	0	0	0	0	49.33	0	
Neutral	39.40	49.38	64.15	40.56	30.06	54.65	50.67	99.91	
Precision	1.00	1.00	1.00	1.00	1.00	1.00	1.00	0.23	0.90
Recall	0.61	0.51	0.36	0.59	0.70	0.45	0.49	1.00	0.59
F1-score	0.75	0.67	0.53	0.75	0.82	0.62	0.66	0.38	0.65
Accuracy	92.28	90.51	88.01	92.07	94.00	89.59	90.29	58.88	86.95

TABLE III. RESULTS OF SVM CLASSIFICATION WITH USING FIRST 1800 PCA COMPONENTS

Emotion	Anger	Contempt	Disgust	Fear	Happy	Sad	Surprise	Neutral	Average
Anger	91.95	0	0	0	0	0	0	0	
Contempt	0	65.63	0	0	0	0	0	0.46	
Disgust	0	0	88.22	0	0	0	0	0.22	
Fear	0	0	0	85.77	0	0	0	0.02	
Happy	0	0	0	0	97.50	0	0	0	
Sad	0	3.75	0	0	0	78.59	0	0.31	
Surprise	0	0	0	0	0	0	94.94	0	
Neutral	8.05	30.63	11.78	14.23	2.50	21.41	5.06	98.99	
Precision	1.00	0.99	1.00	1.00	1.00	0.95	1.00	0.51	0.93
Recall	0.92	0.66	0.88	0.86	0.98	0.79	0.95	0.99	0.88
F1-score	0.96	0.79	0.94	0.92	0.99	0.86	0.97	0.68	0.89
Accuracy	98.87	95.27	98.32	98.01	99.64	96.50	99.28	88.11	96.75

As in Table I Precision values in Table II for all emotions, except Neutral, are close to 1. They are significantly exceed the Recall values, which has values lower than 0.5 in case Disgust, Sad and Surprise emotion states. Due to that F_1 -score is lower than in Table I. For Anger, Fear, Happy F_1 -score is higher than 0.75, which is acceptable. The worst values correspond to Disgust and Neutral emotions states are 0.53 and 0.38, respectively, which is very small. Note that, unlike Table I for all emotions the Precision values is 1, except the Neutral state.

To verify the superiority of AdaBoost over PCA for important feature selection, we increase the count of first PCA components trying get better results. It is increased by about 5 times to 1800 components. This number was selected because it correspond to relatively good accuracy value with small count of features. For bigger number of features the Accuracy does not increase significantly.

Table III shows the results for PCA dimation reduction method with first 1800 principal components. In comparison with Table II Accuracy results in Table III are much better. It exceed 95% in all emotions states, except Neutral, which has 88.11%. This results are comparable with results in Table I. But

F_1 -score is lower, by the reason of small Recall values. The obtained results is much better than for 348 PCA components, but not such good as AdaBoost.

It can be seen that in all Tables the main error of classification caused by the intersection with Neutral emotion state. Our guess is that it caused by absence of some emotion states for the same subject. Frequently, it has Neutral emotion and only one other emotion state.

The results show that the emotion of Contempt has the lowest recognition rate in all three cases. The emotion of Happiness gives the best recognition rate by F_1 -score. The Neutral facial expression state had the lowest F_1 -score value in all three cases, because it is intersect with all cases of emotions. In PCA cases the intersection between the basics emotions is minimal (excluding neutral).

VII. CONCLUSION

We evaluated the accuracy of the proposed algorithms of emotion classification using various evaluation criteria.

Two methods show sufficiently high accuracy: 98% for AdaBoost and 96% for PCA dimensional reduction technics.

Given the obtained values of the accuracy and F_1 -score criteria we can consider these algorithms as state-of-the-art.

The main difference between the algorithms is in use of various approaches to reduce the dimensionality. AdaBoost algorithm does better with dimensionality reduction as it outputs the lesser number of components which provide a high classification accuracy results. Consequently, less time is required to classifier learning and, thus, less time for emotion classification in testing and working stages.

In our future work, we plan to increase the database and improve the quality of classification by using the convolution neural networks, as well as the integration of the system with the analysis of physiological state by video.

ACKNOWLEDGMENT

We would like to thank the reviewers for their comments, and colleagues of the Research Medical and Biological Engineering Center of High Technologies, UrFU.

This work was supported by Act 211 Government of the Russian Federation, contract 02.A03.21.0006 and partially supported by UrFU scholarship for outstanding achievements in science and Russian Foundation for Assistance to Small Innovative Enterprises (FASIE).

REFERENCES

- [1] Peng Yang, Qingshan Liu, Dimitris N. Metaxas, "Boosting encoded dynamic features for facial expression recognition", *Pattern Recognition Letters*, vol.30, no.2, Jan. 2009, pp. 132-139.
- [2] Gwen Littlewort, Jacob Whitehill, Tingfan Wu, Ian Fasel, Mark Frank, Javier Movellan, Marian Bartlett, "The Computer Expression Recognition Toolbox (CERT)", in *Automatic Face & Gesture Recognition and Workshops (FG 2011)*, 2011 IEEE International Conference on. IEEE, Mar. 2011, pp. 298-305.
- [3] Wenfei Gu, Cheng Xiang, Y.V.Venkatesh, Dong Huang ,Hai Lin, "Facial expression recognition using radial encoding of local Gabor features and classifier synthesis", *Pattern recognition*, vol.45, no.1, Jan. 2012, pp. 80-91.
- [4] Patrick Lucey, Jeffrey F. Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, Iain Matthews, "The extended cohn-kanade dataset (CK+): A complete dataset for action unit and emotion-specified expression", *Proceedings of the Third International Workshop on CVPR for Human Communicative Behavior Analysis*, June 2010, pp. 94-101.
- [5] Maja Pantic, Michel Valstar, Ron Rademaker, Ludo Maat, "Web-based database for facial expression analysis", in *International Conference on Multimedia and Expo*, July 2005, pp. 317-321.
- [6] Michael J. Lyons, Shigeru Akemastu, Miyuki Kamachi, Jiro Gyoba, "Coding Facial Expressions with Gabor Wavelets", *3rd IEEE International Conference on Automatic Face and Gesture Recognition*, 1998, pp. 200-205.
- [7] P. Ekman and W. Friesen, *The Facial Action Coding System: A Technique For The Measurement of Facial Movement*. San Francisco: Consulting Psychologists Press, 1978.
- [8] Vahid Kazemi, Josephine Sullivan, "One Millisecond Face Alignment with an Ensemble of Regression Trees", in *Proc. IEEE Conference on Computer Vision and Pattern Recognition*, June 2014, pp. 1867-1874.
- [9] Paul Viola and Michael Jones, "Robust real-time face detection", *International Journal of Computer Vision*, vol.57, no.2, May 2004, pp. 137-154.
- [10] Chang Chih-Chung and Lin Chih-Jen. "LIBSVM: A library for support vector machines", *ACM Transactions on Intelligent Systems and Technology*, vol.2, no.3, Apr. 2011, pp. 1-27.