# Emotion Recognition in the Wild using Deep Neural Networks and Bayesian Classifiers

Luca Surace
University of Calabria - DeMACS
Via Pietro Bucci
Rende (CS), Italy
lucasurace11@gmail.com

Massimiliano Patacchiola
Plymouth University - CRNS
Portland Square PL4 8AA
Plymouth, United Kingdom
massimiliano.patacchiola@plymouth.
ac.uk

Elena Battini Sönmez
Istanbul Bilgi University - DCE
Eski Silahtaraa Elektrik Santral Kazm
Karabekir Cad. No: 2/13 34060 Eyp
Istanbul, Turkey
ebsonmez@bilgi.edu.tr

William Spataro
University of Calabria - DeMACS
Via Pietro Bucci
Rende (CS), Italy
william.spataro@unical.it

Angelo Cangelosi
Plymouth University - CRNS
Portland Square PL4 8AA
Plymouth, United Kingdom
angelo.cangelosi@plymouth.ac.uk

## ABSTRACT

Group emotion recognition in the wild is a challenging problem, due to the unstructured environments in which everyday life pictures are taken. Some of the obstacles for an effective classification are occlusions, variable lighting conditions, and image quality. In this work we present a solution based on a novel combination of deep neural networks and Bayesian classifiers. The neural network works on a bottom-up approach, analyzing emotions expressed by isolated faces. The Bayesian classifier estimates a global emotion integrating top-down features obtained through a scene descriptor. In order to validate the system we tested the framework on the dataset released for the Emotion Recognition in the Wild Challenge 2017. Our method achieved an accuracy of 64.68% on the test set, significantly outperforming the 53.62% competition baseline.

## KEYWORDS

Group emotion recognition; Deep Neural Networks; Bayesian Networks; Ensemble Learning; EmotiW 2017 Challenge

## 1 INTRODUCTION

Automatic emotion recognition has recently become an important research field, due to new possible applications in social media, marketing [24], public safety [5], and human-computer interaction [7]. Emotion recognition is generally achieved through the analysis of facial muscles movements, often called action units. After isolating the face of the subject, it is possible to assign it an emotion using action units analysis. Despite noteworthy results in structured condition this approach becomes unfeasible in unstructured environments where multiple factors (e.g. occlusions, variable lighting conditions, image quality, etc.) may affect recognition.

During the years there have been many attempts to build robust methods. For example, in [29] the authors used histogram of gradients in order to address the problem of human emotion identification from still pictures taken in semi-controlled environments. In [25] the influence of multiple factors (pose, resolution, global and local features) on different facial expressions was investigated. The authors used an appearance based approach dividing the images into sub-blocks and then used support vector machines to learn

pose dependent facial expressions. Deep neural networks have been used in [26]. The authors started from a network pre-trained on the generic ImageNet dataset, and performing supervised fine-tuning in a two-stage process. This cascading fine-tuning achieved better results compared to a single stage fine-tuning. A significant contribution to the use of deep neural networks for emotion recognition was given in [22] and [36]. The authors demonstrated the strength of this approach achieving state-of-art performances on the CK+ dataset [23]. The classification of the emotion of a group of people is a different task. Previous research [9] focused on the development of two parallel approaches for measuring happiness level: top-down and bottom-up. One of the first articles which considered the structure of the scene as a whole is presented in [14]. The authors showed that the structure of the group provides meaningful context for reasoning about the individuals and they considered the group structure from both the local and the global point of view. A complete review of all the methods is out of the scope and we refer the reader to a recent survey [33].

Taking into account past literature we propose a method which integrates both global and local information. A bottom-up module isolates the human faces which are present in the picture and gives them as input to a pre-trained Convolutional Neural Network (CNN). At the same time a top-down module finds the label associated with the scene and passes them to a Bayesian Network (BN) which estimates the posterior probabilities of each class. The output of the system is the probability of the image to belongs to three different classes: positive, neutral, and negative. Experiments were conducted on different architectures, achieving the best results with a pipeline that redirects the output of the CNN to the Bayesian classifier. We tested the system on the dataset released for the Emotion Recognition in the Wild Challenge 2017 (EmotiW) [8, 10, 11], obtaining an accuracy of 67.75% on the validation set, and 64.68% on the test, outperforming the baseline of the competition [11].

## 2 PROPOSED METHOD

Our method is based on the idea that the group emotion can be inferred using both top-down [3] and bottom-up [10] approaches. The former considers the scene context, such as background, clothes,

location, etc. The latter estimates the face expressions of each person in the group. We can summarize the pipeline for the bottom-up module in three steps:

(1) Face detection
(2) Face pre-processing
(3) CNN forward pass

The first step is the face detection, which has been obtained through a commercial library [1]. This step returns a list of frames containing isolated human faces. In the second step the faces are cropped, scaled and normalized. Finally in the third step the cropped faces are given as input to a pre-trained CNN and the output for each class are estimated through a forward pass. In parallel it is possible to run the top-down module which consists of three steps:

(1) Acquiring the scene descriptors
(2) Set evidences in the BN
(3) Estimate the posterior distribution of the BN

In the first step the top-down module estimates the content of an image returning descriptors which may be context-specific (e.g. party, protest, festival, etc.) and group-specific (e.g. team, military, institution, etc.). This step is achieved through the same previously adopted library, but it can be easily replaced with a state-of-the-art algorithm [40]. In the second step, the descriptors are considered as observed variables and the corresponding nodes in the Bayesian network are set as evidence. In the third and last step the posterior distribution of the root node in the BN is estimated using the belief propagation algorithm [32]. In the next sections, details of the proposed methodology are reported, together with an overview of the entire system (Figure 1).
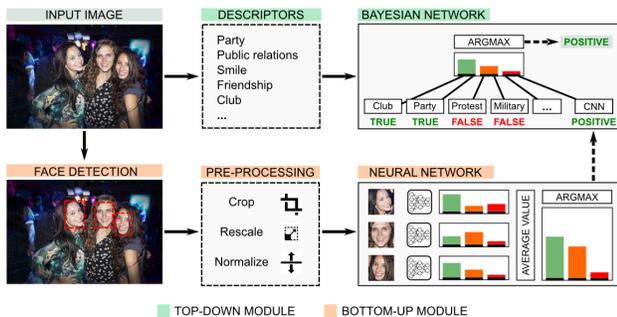


**Figure 1: Overview of the proposed system. In the top-down module (light-green) the scene descriptor returns a list of labels used by the BN. In the bottom-up module (light-orange) the face detector isolates the faces which are given as input to the deep neural network after a pre-processing phase.**

## 2.1 Bottom-up module

The bottom-up module uses a CNN to estimate the emotion of isolated human faces. Deep neural networks recently achieved outstanding performances in a variety of tasks such as speech recognition [17, 18], and head pose estimation [31]. We trained several networks having multiple architectures. The best results

---

[1]Google Vision API

have been achieved with a variant of AlexNet [21]. The input to the network is a color image of size 64×64 pixels which is passed through 7 layers: 3 convolutional layers, 2 sub-sampling layers, 1 fully connected layer and the final output layer. The first convolution layer produces 64 feature maps, applying a convolutions kernel of size $11 \times 11$. The second convolution layer generates 128 feature maps, using a convolutional kernel of size $5 \times 5$. The third layer produces 256 feature maps via a convolutional kernel of size $3 \times 3$. The sub-sampling layers use max-pooling (kernel size $3 \times 3$) to reduce the image in half. The result of the third convolution is given as input to a fully connected layer (512 units). Finally, the network has three output units representing the three emotions: positive, neutral, negative. The first two convolutional layers are normalized with local response normalization [21]. A rectified linear unit activation function is applied to each convolutional layer and to the first fully connected layer. For the training we used an adaptive gradient method, RMSProp [38], and the balanced batches technique proposed in [36]. As a loss function we used the softmax cross entropy [35] between the target $t$ and the estimated value $o$, defined as follows:

$$J(T, O) = -\frac{1}{N} \sum_{n=1}^{N} \left[ t_n \ln(o_n) + (1 - t_n)\ln(1 - o_n) \right] \quad (1)$$

where $N$ is the size of the batch, $T = \{t_1, ..., t_N\}$ is the set of target values, $O = \{o_1, ..., o_N\}$ is the set of output values. Once the network has been trained it is possible to estimate the average group emotion for the faces present in the image. First of all, we averaged the predictions resulting from a forward pass on all the input faces, similarly to [39]. Secondly, we returned the class corresponding to the higher value. We can summarize the two steps in a single equation:

$$\hat{o} = \text{argmax} \left( \frac{\sum_{k=1}^{K} \sigma(\mathbf{o}_k)}{K} \right) \quad (2)$$

where $K$ is the total number of faces, $\mathbf{o}$ is a three dimensional vector representing the output of the network, and $\sigma$ is the softmax function. The resulting scalar $\hat{o}$ represents the index of the class which better represents the scene emotion, based on all the faces that has been found in the image. This method can be extremely effective, but is has some drawbacks. First of all, it requires to identify at least one face per image. Secondly, the faces should be a good predictor of the overall emotion, which is not always the case. To compensate this source of error we used a top-down module which helps to describe the scene.

## 2.2 Top-down module

Previous literature shows that global scene information is very useful for group emotion classification [10]. In particular, scenarios focusing only on small details can easily lead to miss-classification errors. This is why in the top-down module we used a whole-scene descriptor. In this section we describe the procedure for collecting the descriptors and how they have been integrated in the BN. In a preliminary phase we collected meaningful descriptors for each image contained in the training set. Subsequently, we built an histogram of descriptors for the three emotions. We found a total

of 812 scene descriptors, appearing with a certain frequency in the dataset. The descriptors are represented in a word cloud in Figure 2. It is possible to see how descriptors such as smile, friendship, festival, etc. recur with an high frequency in positive images and less frequently in the neutral and negative groups.



**Figure 2: Overview of the descriptors for each one of the three classes: positive (left-green), neutral (center-orange), negative (right-red). Words which have a larger font size, recur more frequently.**

The descriptors obtained through the preliminary phase have been used as dependent nodes in a BN. BNs represent a valid formalism to model probabilistic relationships between several random variables and their conditional dependencies via a directed acyclic graph. They have been used in different applications such as medical diagnosis [28], and cognitive modeling [30] (for a review see [13]). In this work we started from the assumption of independence between each pair of descriptors. This assumption is an oversimplification because does not capture relationships between different features. However, from a practical point of view it works extremely well, and it is applied in many state-of-the-art classifiers. From the mathematical point of view we want to estimate $\hat{y}$, the class having the higher probability given the observed descriptors:

$$\hat{y} = \underset{y}{\operatorname{argmax}} \, P(y|x_1, x_2, ..., x_N) \qquad (3)$$

The posterior distribution associated with the root node $y$ is proportional to the product of the prior $P(y)$ and the likelihood $P(x_i|y)$ of each one of the $N$ dependent variables $x$, as follows:

$$P(y|x_1, x_2, ..., x_N) = P(y) \prod_{i=1}^{N} P(x_i|y) \qquad (4)$$

In our particular case the root node is a multinomial probability distribution which models the three possible outcomes: positive, neutral and negative. For each descriptor there is a dependent variable, which is modeled with a Bernoulli distribution: true (present), false (absent). The main point for obtaining $\hat{y}$ is to find the conditional probabilities $P(x_i|y)$. Since for the training set we know the emotion associated with every image, we can use Maximum Likelihood Estimation (MLE) [34] to find the conditional probability distributions for each node. For instance, naming $N_t^+$ the number of time a specific descriptor $x_i$ has been counted in conjunction with positive images, and $N_f^+$ the number of time that descriptor has not been counted, we can estimate the conditional probability as follows:

$$P(x_i = \text{true}|y = \text{positive}) = \frac{N_t^+}{N_t^+ + N_f^+}$$
$$P(x_i = \text{false}|y = \text{positive}) = \frac{N_f^+}{N_t^+ + N_f^+} \qquad (5)$$

For the same descriptor we can also estimate the probability $P(x_i|y = \text{neutral})$ of being associated with a neutral emotion, and $P(x_i|y = \text{negative})$ the probability of being associated with a negative emotion. Using these probabilities we can estimate the Bernoulli distributions associated with every descriptor and build the corresponding conditional probability tables in the BN.

## 2.3 Integration

There are different ways the results of the two modules can be combined. For instance, considering the two modules as a committee of experts we can use ensemble averaging to reduce the error of the models. Another possibility is to redirect the value obtained by the bottom-up module as input to the BN in the top layer. After a preliminary research we decided to adopt the second solution.

Going back to Equation 3 and 4, we can hypothesize the presence of an additional input feature $x_{N+1}$ having as prior a three-categorical multinomial distribution (positive, neutral, negative). In order to integrate the new node in the BN it is necessary to estimate the conditional probability table for $P(x_{N+1}|y)$. Similarly to Equation 5 it is possible to use MLE to find the conditional distributions for the dependent node. In the particular case considered here the conditional distribution is represented by the confusion matrix obtained testing the network on the dataset.

## 3 EXPERIMENTS

In this section we describe the methodology followed during the training phase and the results achieved. We evaluate the proposed method on the GAF database [10]. This database consists of 6470 total images, which have been divided in 3633 images for the training set, 2065 for the validation set, and 772 for the test set. The dataset contains images obtained from social networks captured during social events. These images can be from positive social events (marriages, parties, etc.), neutral event (meetings, convocations, etc.), or negative events (funeral, protests, etc). The baseline score for this dataset has been obtained using the CENTRIST [41] approach and support vector regression. CENTRIST is a scene descriptor and is computed on the whole image. It takes into consideration both the bottom-up and top-down attributes. The classification accuracy is used as the metric in the challenge. The model baseline achieved 52.97% on the Validation set and 53.62% on the test set.

## 3.1 Methods

The CNNs have been trained on the dataset available with the challenge. For each one of the images in the training set we isolated the faces and performed different pre-processing operations. First of all the faces have been cropped and re-scaled to $64 \times 64$ pixels. Then a min-max normalization has been applied. During the training we randomly selected a balanced batch of 63 images (21 for each emotion) and performed gradient descent for 1500 epochs. As optimizer we use the RMSProp [38], which has been selected after

a preliminary comparison between other methods such as Adagrad [12] and Adam [20]. A learning rate $\alpha = 10^{-3}$, a decay of 0.9, and $\epsilon = 10^{-10}$ were adopted. The weights of the CNNs have been initialized using the Xavier initialization method [16]. We used dropout [37] with 0.5 probability in between the internal layers in order to prevent overfitting.

We implemented the algorithm in Python using the Tensorflow library [1] for the CNNs training, and OpenCV [19] for the pre-processing operations on the images. Experiments were carried out on a workstation having 16 cores processor, 32 GB of RAM, and the NVIDIA Tesla K40 graphical processing unit. On this hardware the training of the CNN took approximately 8 minutes. The evaluation on the validation set (2065 images) using the whole pipeline took 325 minutes.

## 3.2 Results

We obtained the best performance with the ensemble method, which lead to an accuracy of 67.75% on the validation set and 64.68% on the test set. Those results are significantly higher than the challenge baseline accuracy of 52.97% (validation) and 53.62% (test). Comparative results between the BN-only, CNN-only and ensemble approaches are showed in Figures 3. The ensemble method outperformed the results obtained using the isolated modules, supporting previous work that demonstrated how an ensemble can improve the performance of emotion recognition systems [15].

**Figure 3: Validation accuracy comparison for the stand-alone solutions and the complete system (BN and CNN).**

Considering the confusion matrices for the three methods, reported in figure 4, we can see how the integration of BN and CNN leads to the best performance. The darker cells are on the main diagonal, meaning that the system can associate unknown input features to correct labels.

## 4 CONCLUSIONS AND FUTURE WORK

In this article we investigated the use of deep CNNs and Bayesian classifiers for group emotion recognition in the wild. Our method uses an approach which takes into account both top-down and bottom-up methods. The top-down module estimates the group emotion based on scene descriptors, which are integrated in a BN. On the other hand the bottom-up module identifies human faces in the picture and returns an average emotion estimation. The output of the bottom-up module is then redirected to the BN in the higher

**Figure 4: Confusion matrices for (a) BN on validation set; (b) CNN on validation set; (c) complete system on validation set; (d) complete system on test set. Color scale is based on the accuracy value.**

layer and considered as a dependent node. The method has been tested on the EmotiW'17 challenge dataset, obtaining an accuracy of 67.75% on the validation set and 64.68% on the test set, and achieving a significant improvement over the baseline performance [10].

Future work should focus on different aspects which may have an important role on the accuracy. A more sophisticated approach for integrating the output of the CNN for each detected face should be taken into account. For example in [9] and [39], a weighted average based on the size of the face is used in order to calculate an overall score. This method is reasonable since it gives a lower weight to faces which are on the background, and which may carry less information. In order to further improve the classification accuracy obtained through the CNN it is possible to use a divide-and-conquer strategy. Instead of relying on a single CNN to estimate the three categories, it may be possible to split the classifier in three sub-networks which are specialized in the identification of a single emotion. Such an approach showed major improvements in different tasks [4, 6, 27, 31].

Another critical point is the integration of the estimates made by the two modules. In this work we empirically found that redirecting the output of the bottom-up module to the BN in the top layer leads to a slight improvement. However, other methods, for instance bagging [2], should be considered in this delicate phase. In conclusion, further research is needed in order to understand which approach may lead to higher performances.

## REFERENCES

[1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dan Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike

Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. 2015. TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems. (2015). http://tensorflow.org/ Software available from tensorflow.org.

[2] Leo Breiman. 1996. Bagging predictors. *Machine learning* 24, 2 (1996), 123–140.

[3] Aleksandra Cerekovic. 2016. A deep look into group happiness prediction from images. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction.* ACM, 437–444.

[4] Sung-Bae Cho and Jin H Kim. 1995. Combining multiple neural networks by fuzzy integral for robust classification. *IEEE Transactions on Systems, Man, and Cybernetics* 25, 2 (1995), 380–384.

[5] Chloé Clavel, Ioana Vasilescu, Laurence Devillers, Gaël Richard, and Thibaut Ehrette. 2008. Fear-type emotion recognition for future audio-based surveillance systems. *Speech Communication* 50, 6 (2008), 487–503.

[6] Edward Collins, Sushmito Ghosh, and Christopher Scofield. 1988. An application of a multiple neural network learning system to emulation of mortgage under-writing judgments. In *Proceedings of the IEEE International Conference on Neural Networks*, Vol. 2. 459–466.

[7] Roddy Cowie, Ellen Douglas-Cowie, Nicolas Tsapatsoulis, George Votsis, Stefanos Kollias, Winfried Fellenz, and John G Taylor. 2001. Emotion recognition in human-computer interaction. *IEEE Signal processing magazine* 18, 1 (2001), 32–80.

[8] Abhinav Dhall et al. 2012. Collecting large, richly annotated facial-expression databases from movies. (2012).

[9] Abhinav Dhall, Roland Goecke, and Tom Gedeon. 2015. Automatic group happiness intensity analysis. *IEEE Transactions on Affective Computing* 6, 1 (2015), 13–26.

[10] Abhinav Dhall, Jyoti Joshi, Karan Sikka, Roland Goecke, and Nicu Sebe. 2015. The more the merrier: Analysing the affect of a group of people in images. In *Automatic Face and Gesture Recognition (FG), 2015 11th IEEE International Conference and Workshops on*, Vol. 1. IEEE, 1–8.

[11] Abhinav Dhall, Jyoti Joshi, Karan Sikka, Roland Goecke, and Nicu Sebe. 2017. From Individual to Group-level Emotion Recognition: EmotiW 5.0, ACM ICMI 2017, Vol. in press.

[12] John Duchi, Elad Hazan, and Yoram Singer. 2011. Adaptive subgradient methods for online learning and stochastic optimization. *Journal of Machine Learning Research* 12, Jul (2011), 2121–2159.

[13] Nir Friedman, Dan Geiger, and Moises Goldszmidt. 1997. Bayesian Network Classifiers. *Machine Learning* 29, 2 (01 Nov 1997), 131–163. https://doi.org/10.1023/A:1007465528199

[14] Andrew C Gallagher and Tsuhan Chen. 2009. Understanding images of groups of people. In *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on.* IEEE, 256–263.

[15] Michael Glodek, Stephan Tschechne, Georg Layher, Martin Schels, Tobias Brosch, Stefan Scherer, Markus Kächele, Miriam Schmidt, Heiko Neumann, Günther Palm, et al. 2011. Multiple classifier systems for the classification of audio-visual emotional states. *Affective Computing and Intelligent Interaction* (2011), 359–368.

[16] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics.* 249–256.

[17] Alex Graves, Abdel-rahman Mohamed, and Geoffrey Hinton. 2013. Speech recognition with deep recurrent neural networks. In *Acoustics, speech and signal processing (icassp), 2013 ieee international conference on.* IEEE, 6645–6649.

[18] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdel-rahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, et al. 2012. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine* 29, 6 (2012), 82–97.

[19] Itseez. 2015. Open Source Computer Vision Library. https://github.com/itseez/opencv. (2015).

[20] Diederik Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[21] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems 25*, F. Pereira, C. J. C. Burges, L. Bottou, and K. Q. Weinberger (Eds.). Curran Associates, Inc., 1097–1105. http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

[22] André Teixeira Lopes, Edilson de Aguiar, Alberto F De Souza, and Thiago Oliveira-Santos. 2017. Facial expression recognition with Convolutional Neural Networks: Coping with few data and the training sample order. *Pattern Recognition* 61 (2017), 610–628.

[23] Patrick Lucey, Jeffrey F Cohn, Takeo Kanade, Jason Saragih, Zara Ambadar, and Iain Matthews. 2010. The extended cohn-kanade dataset (ck+): A complete dataset for action unit and emotion-specified expression. In *Computer Vision and Pattern Recognition Workshops (CVPRW), 2010 IEEE Computer Society Conference on.* IEEE, 94–101.

[24] Daniel McDuff, Rana El Kaliouby, Thibaud Senechal, David Demirdjian, and Rosalind Picard. 2014. Automatic measurement of ad preferences from facial responses gathered over the internet. *Image and Vision Computing* 32, 10 (2014), 630–640.

[25] S Moore and R Bowden. 2011. Local binary patterns for multi-view facial expression recognition. *Computer Vision and Image Understanding* 115, 4 (2011), 541–558.

[26] Hong-Wei Ng, Viet Dung Nguyen, Vassilios Vonikakis, and Stefan Winkler. 2015. Deep learning for emotion recognition on small datasets using transfer learning. In *Proceedings of the 2015 ACM on international conference on multimodal interaction.* ACM, 443–449.

[27] Hanh H Nguyen and Christine W Chan. 2004. Multiple neural networks for a long term time series forecast. *Neural Computing & Applications* 13, 1 (2004), 90–98.

[28] Agnieszka Onisko, Marek J Druzdzel, Hanna Wasyluk, et al. 1999. A Bayesian network model for diagnosis of liver disorders. In *Proceedings of the Eleventh Conference on Biocybernetics and Biomedical Engineering*, Vol. 2. 842–846.

[29] Carlos Orrite, Andrés Gañán, and Grégory Rogez. 2009. Hog-based decision tree for facial expression classification. *Pattern Recognition and Image Analysis* (2009), 176–183.

[30] Massimiliano Patacchiola and Angelo Cangelosi. 2016. A developmental Bayesian model of trust in artificial cognitive systems. In *Development and Learning and Epigenetic Robotics (ICDL-EpiRob), 2016 Joint IEEE International Conference on.* IEEE, 117–123.

[31] Massimiliano Patacchiola and Angelo Cangelosi. 2017. Head Pose Estimation in the Wild using Convolutional Neural Networks and Adaptive Gradient Methods. *Pattern Recognition* (2017).

[32] Judea Pearl. 2014. *Probabilistic reasoning in intelligent systems: networks of plausible inference.* Morgan Kaufmann.

[33] Evangelos Sariyanidi, Hatice Gunes, and Andrea Cavallaro. 2015. Automatic analysis of facial affect: A survey of registration, representation, and recognition. *IEEE transactions on pattern analysis and machine intelligence* 37, 6 (2015), 1113–1133.

[34] FW Scholz. 1985. Maximum likelihood estimation. *Encyclopedia of statistical sciences* (1985).

[35] John Shore and Rodney Johnson. 1981. Properties of cross-entropy minimization. *IEEE Transactions on Information Theory* 27, 4 (1981), 472–482.

[36] Elena Battini Sönmez and Angelo Cangelosi. 2017. Convolutional neural networks with balanced batches for facial expressions recognition. In *Ninth International Conference on Machine Vision.* International Society for Optics and Photonics, 103410J–103410J.

[37] Nitish Srivastava, Geoffrey E Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. 2014. Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research* 15, 1 (2014), 1929–1958.

[38] Tijmen Tieleman and Geoffrey Hinton. 2012. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning* 4, 2 (2012).

[39] Vassilios Vonikakis, Yasin Yazici, Viet Dung Nguyen, and Stefan Winkler. 2016. Group happiness assessment using geometric features and dataset balancing. In *Proceedings of the 18th ACM International Conference on Multimodal Interaction.* ACM, 479–486.

[40] Yilin Wang, Suhang Wang, Jiliang Tang, Huan Liu, and Baoxin Li. 2015. Unsupervised Sentiment Analysis for Social Media Images.. In *IJCAI.* 2378–2379.

[41] Jianxin Wu and Jim M Rehg. 2011. CENTRIST: A visual descriptor for scene categorization. *IEEE transactions on pattern analysis and machine intelligence* 33, 8 (2011), 1489–1501.