

LITERATURE REVIEW

By

Aaron Ward - B00079288

Supervisor(s): Stephen Sheridan

SUBMITTED IN PARTIAL FULFILLMENT OF THE
REQUIREMENTS FOR THE DEGREE OF
B.SC IN COMPUTING AND INFORMATION TECHNOLOGY
AT
INSTITUTE OF TECHNOLOGY BLANCHARDSTOWN
DUBLIN, IRELAND
2017

Declaration

I hereby certify that this material, which I now submit for assessment on the programme of study leading to the award of **B.Sc in Computing and Information Technology** in the Institute of Technology Blanchardstown, is entirely my own work except where otherwise stated, and has not been submitted for assessment for an academic purpose at this or any other academic institution other than in partial fulfillment of the requirements of that stated above.

Dated: 2017

Author:

Aaron Ward - B00079288

Introduction

The following literature review will provide and insight and analysis into the current research regarding machine learning and artificial intelligence, specifically on the topic of neural networks. The paper shall consist of a review of five papers based on computer vision, convolutional neural networks or sentiment analysis research. This is followed by a related work section, which shall display the relevance of these papers to the proposed research given and the similarities between the two. Furthermore, a critical analysis shall be given into the quality of these papers, outlining some inconsistencies and weak points, if any. This paper will not explore the papers based on the theory of neural networks, but the research that implements them.

Literature

Analyzing and Detecting Employee's Emotion for Amelioration of Organizations

Subhashini and Niveditha make the opening statement that emotions usually do not any place in a work environment in current society. Although the expression of feelings is suppressed in places of work, they suggest that emotions can affect five major areas in competitive advantage. The five given aspects of competitive advantage are as follows: Intellectual Capital, Customer Service, Organizational Reactivity, Production, Employee appeal and retentivity [Subhashini and Niveditha, 2015]. In order to counter this apprehensiveness to expression of emotions in the work place, Subhashini and Niveditha suggest the concept of a facial emotion tracking system that will map the facial expressions of an employees face as they enter the organization. Two related works are also given; Emotion Detections Based on Text and Emotion Recognition Based on Brain-Computer Interface Systems.

The system architecture given by Subhashini and Niveditha briefly describes the program. Most employees entering a building to an organization must swipe a card to clock into the work hours. They suggest that they have designed a new system that removes the need for card swiping but also implements emotion detection. The employee must look into a camera that will prove their presence in the vicinity but will also perform some emotion detection. The system was implemented in the C Sharp programming language and uses skin tone segmentation to detect. The binary image is then converted to an RGB image and an inspection of every individual pixel if performed. If the RGB value is greater than 110, then the pixel colour is refactored to be a white pixel, other wise it becomes a black pixel. This is done to make it easy to detect facial features in the video stream [Subhashini and

Niveditha, 2015]. Once detected, the image around the face is cropped. They then apply a Bezier Curve to the regions around the lips and eyes of the person being analyzed. The results of the persons identity and emotional status are then stored within a database. They conclude their paper by explaining that this system can be used by management to gain an understanding of their employees sentimental state.

Deep Learning for Video Classification and Captioning

Wu et al.'s paper provides an in depth analysis on the methods for video classification and video captioning in terms of deep learning. They claim that because of the exponential growth in internet bandwidth and computing power, video communications are becoming more and more prevalent, therefore paving the way for new video understanding applications [Wu et al., 2016]. They make reference to current implementations to prove the growth of interest in the field of computer vision and video analysis, notably the ImageNet challenge.

Move over, they go on to give brief description of the two "deep learning modules" that have been used for visual analysis: Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN). It is explained that LeCun et al. that developed LeNet-5 made a break through when they developed a CNN using the Back-Propogation Algorithm. But its noted that this is limited in performance when the complexity of the tasks is increased. Deep belief networks were developed to train networks in a unsupervised manner in order to counter this problem [Wu et al., 2016]. AlexNet, a CNN proposed by Krizhevsky et al. in 2012, introduced two way to increase the performance of CNN's using ReLU (Rectified Linear Units) and Dropout to descrease overfitting [Wu et al., 2016]. Seondly, RNN's are brought forward. Wu et al. explain the difference between CNN's and RNN's, stating that CNN's are all feed forward networks that do not use cycling, which can prove be disadvantageous when working with sequence labeling. Two issues can occur with RNN's: Vanishing Gradient and Exploding Gradients as short term memory is used when cycling through the network. The solution given to this is an RNN variant called Long Short-Term Memory (LSTM).

Image-Based Video Classification using CNN's and LSTM's

They state that Karparthy et al. researched the common architectures for learning spatial-temporal clues in large video datasets. It appeared that models using single frames as input achieve similar results as models using stacks of frames. From this, Simonyan and Zisserman proposed the idea of the Two-Stream Approach, because of the cost effectiveness and time consumption that come with training 3D CNN's. This Two-Stream approach involves training the CNN on single and stack frames concurrently. Both outputs are put through a score fusion. The result is the weighted sum of both scores [Wu et al., 2016].

Although Two-Stream is a good approach, it is not sufficient as it is not capable of dealing with long video clips. Therefore, LSTM's are utilized as they do not suffer from the problem of vanishing gradients. It has been found that CNN's and LSTM's compliment each other when working in conjunction with each other [Wu et al., 2016]. They conclude their paper with a summary of the written topics about, regarding the growth for the need for video understanding applications, the used of CNN's and RNN's, in addition to the variants of these deep learning modules.

Subject independent facial expression recognition with robust face detection using a convolutional neural network

As stated by Matsugu et al., difficulties may arise with facial recognition. In terms of a face being in a smiling-like state, could have different implications. As well as this, a facial recognition system should be able to work with a wide range of variability of faces. They then give some examples of past implementations such as facial recognition with rigid head movement by Black and Yackoob in 1995, and speak about how this does not meet the requirements of dealing with wide variance. A rule based system is proposed [Matsugu et al., 2003]. With their model, layer trains on a module-by-module (module being the nose, eyes, mouth etc) basis. Meaning each layer trains on a certain facial feature. Each of the neurons perform an averaging of some local receptive fields then they use a skin tone detector to detect each module on the face [Matsugu et al., 2003]. For training, Layer one and layer two are trained for 8 modules using back propagation. Layer three and four train on more

complex feature detectors such as the mouth and eyes. The output is then sent to the rule based algorithm for handling variability and robustness.

The rule based algorithm takes the output of the CNN and measures the distance between the features. From these calculation, the rules are applied to determine if the person is in a laughing or smiling state. The rules are summarised as follows [Matsugu et al., 2003]:

- The distance between eyes and lip get shorter.
- The horizontal length of lip gets longer.
- The eyes wrinkle.
- The gradient of lip from the end point to the end point increases.
- Detection of teeth increases.
- The edges (wrinkles) of cheek increase.

In conclusion, they received a 97.6 percent accuracy for 10 test subjects with 5600 images. They assert that their model is significantly more efficient as they only require one CNN due to their rule based algorithm, in contrast to Fasels implementation that uses two CNN's working in synergy with one another.

Neuromarketing - The Art and Science of Marketing and Neurosciences Enabled by IoT Technologies

A paper by Arthmann and Li describes the growing field of Neuromarketing with an opening statement: Advertisers recognise that there is a relationship between stimulating the emotions of a customer and influencing their actions. Online shopping has drastically affected the store sales, and it is also explained that more than 3500 stores have shut down due to bankruptcy in 2017 [Arthmann and Li, 2017]. Their answer to this change in buying is Neuro Linguistic Programming. Neuro Linguistic Programming (NLP), not to be confused with natural language processing, is a form of observing the verbal and non verbal communication of humans. Eye accessing cues (eye movements) are said to be linked to certain

emotions or thoughts. Neuromarketing incorporates NLP and IoT devices to understand the consumer sentiment more extensively. Neuromarketing aims to remove marketing biases by utilizing the consumers subconscious. One example given by Arthmann and Li is the notion of facial coding and motion tracking, which is put in place to determine why consumer make certain decisions. Furthermore, they make their belief clear of retailers benefiting from this when it is put forward that artificial intelligence and machine learning will evolve, giving better results of consumer preference shifts. Additionally, These neuromarketing systems can replace thermoimaging people counting devices that clock people walking into stores that may not be eligible customers (children). Further examples are provided for these technologies. The describe a scenario when a customer is given an image of a product and a price in front of a webcam, and by performing facial coding and sentiment analysis, we may get a better sense of what the consumer is feeling. Arthmann and Li conclude their paper by declaring the future potential of these systems using when integrated with "always on" IoT technologies.

Relevant Work and Critical Analysis

Subhashini and Niveditha, 2015's work on analyzing and detecting employees emotions for organizations ties in very well with the proposed project as it involves using sentiment analysis on subject to gain an underlining understanding of their emotions that may or may not be expressed verbally. Their approach to skin tone segmentation for detecting a humans face may prove very beneficial to the proposed project. The main research goal was to achieve employee identification in conjunction with facial emotional analysis and they were successful in execution, They used a Bezier Curve on the subjects lips and eye to detect the emotion expressed by analysing the gradient of the curve. However, Although this paper proved that the project was a success, there are some inconsistencies. For example: in the related work, they do not explain how the work is related and only give titles. Secondly, the results show no code snippets or pseudocode to explain the implementation of the system. Only screen shots of the user interface are provided. Furthermore, there are more absences of proof. The paper is lacking statistics and graphs to display the accuracy or progress of the systems performance and there are some bold statements used that are not back up by citations like when it is said that "emotions were considered a forbidden topic in the working place". Despite these weak areas in the paper, a good aspect of this system is the use of real life subjects used in testing.

The review of Deep Learning for Video Classification and Captioning by Wu et al. provides an in depth look into the aspects of different neural networks and what are they strong and weak points. The motivation for their research is driven by their claim that video communications is growing and that their needs to be better applications for video understanding. The relevancy of this paper provides the concept of the "Two-stream" architecture. Although it is not planned to develop two convolutional neural networks (CNN), is it sought after to

develop a score fusion algorithm, similar to the one mentioned in this paper. A CNN will be developed and the application will utilize a tone analyser for voice sentiment analysis and the two scores will be combined by such an algorithm to provide the weighted sum of the two scores. As this paper is very in depths and draws good comparisons between the different techniques that can be used for video classification. Despite the quality of this paper, some aspects need improvement. There is heavy usage of words like "we" used. Also, some statements are made by [Wu et al., 2016] that are not cited to supported their claim. This is evident when its said "As deep learning for video analysis is an emerging and vibrant feild..." (According to whom?).

Subject independent facial expression recognition with robust face detection using a convolutional neural network by Matsugu et al. illustrates the difficulties that may arise when performing facial recognition. They highlight the problems that may occur in terms of being able to handle variability of subject faces, and certain angles. Their approach to this problem is addressed by implementing a rule based algorithm that analysis the results given from the CNN. Furthermore, their model is designed to be segment and be trained on specific facial features instead the face as a whole. Their model proves to be a success as they score an accuracy of 97.6 percent, also they do not require a second CNN working concurrently to achieve similar results as other models have done so previous, which can be cost effective. Some similarities arise between this paper and the proposed project, they both use sentiment analysis and require the ability to handle a wide range of variability. This proves beneficial to the proposed project as it provides inspiration to used a rule based algorithm for determining emotions. Even though this paper is well written, there are some issues. In certain parts there are abbreviations to words given without the full word be given prior which can cause confusion to the reader. For example: "FP neurons". In addition to this, their model is specific for smiling faces and doesn't accommodate for other emotions, which should be at least provided in a further work section.

Arthmann and Li's paper titled Neuromarketing The Art and Science of Marketing and Neurosciences Enabled by IoT Technologies is a promising insight to the field of neuromarketing. They recognize the association of online shopping and loss of sale for retail stores, and give example of how neuro linguistic programming and IoT technologies can be used as a combination to understand their customers sentiment. Furthermore, it's heavily argued

that the integration of AI and machine learning will evolve this concept to understand the thoughts of consumers and further tackle loss in productivity. This proves relevant to the proposed project as it is very similar. They both aim to gain a deeper understanding of human sentiment that may not be express verbally. Although this is possibly the most interesting paper of this review, it is severally lacking citations. Also, there are a lot of assumption brought forward with no clear indication as to how this knowledge is known. Additionally, it is also assumed that people will adopt these "always on" IoT devices and agree for their physical aspects to be used for consumer targeted marketing.

Conclusion

In conclusion, this paper has reviewed five papers in relation to convolutional neural networks, facial sentiment analysis, emotion detection and current applications in the real world. Relevant aspects include detecting employees emotions in a work environment, implementing a score fusion algorithm for achieve a summation of two sentiment detecting technologies, the idea of segmenting facial features while training to accommodate for variability and the notion using human emotion understanding for a business solution.

Bibliography

- Christopher Arthmann and I-Ping Li. Neuromarketing the art and science of marketing and neurosciences enabled by iot technologies, 2017. URL https://www.iiconsortium.org/pdf/2017_JoI_Neuromarketing_IoT_Technologies.pdf.
- Masakazu Matsugu, Katsuhiko Mori, Yusuke Mitari, and Yuji Kaneda. Subject independent facial expression recognition with robust face detection using a convolutional neural network. *Neural Networks*, 16(5):555 – 559, 2003. ISSN 0893-6080. doi: [https://doi.org/10.1016/S0893-6080\(03\)00115-1](https://doi.org/10.1016/S0893-6080(03)00115-1). URL <http://www.sciencedirect.com/science/article/pii/S0893608003001151>. Advances in Neural Networks Research: IJCNN '03.
- R. Subhashini and P.R. Niveditha. Analyzing and detecting employee's emotion for amelioration of organizations. *Procedia Computer Science*, 48(Supplement C):530 – 536, 2015. ISSN 1877-0509. doi: <https://doi.org/10.1016/j.procs.2015.04.131>. URL <http://www.sciencedirect.com/science/article/pii/S1877050915006407>. International Conference on Computer, Communication and Convergence (ICCC 2015).
- Zuxuan Wu, Ting Yao, Yanwei Fu, and Yu-Gang Jiang. Deep learning for video classification and captioning. *CoRR*, abs/1609.06782, 2016. URL <http://arxiv.org/abs/1609.06782>.