

Detection of cats in images

ANNA MARIA WALACH (s121540)

Technical University of Denmark

02238 Biometric Systems

More and more attention is focused on feral and wild animals, providing them with medical care, food, checking their safety. It is possible with live streams and cams, that allows to increase awareness of thousands of people, but also make life of people taking care of those animals significantly easier, by being able to constantly monitor the status of pupils. Possibility of automatic detection and recognition of animals in such video signals could significantly increase ability to reliable counting of populations or checking for medical care needs. This research is performing comparison of current methods in the field and applies open source algorithm to simple date to estimate the robustness of out-of-the-box implementation.

Categories and Subject Descriptors: I.5.4 [**Pattern Recognition**]: Applications—*computer vision*; H.4.0 [**Information Systems Applications**]: General

General Terms: Algorithms, Experimentation

Additional Key Words and Phrases: cat, detection, images

1. INTRODUCTION

There are more and more projects that uses live streaming and social media power to help and safe feral cats. Some of them are meant to control the population in the area [LangleyAnimalProtectionSociety 2015], other focus on raising awareness about feral cats, spaying and neutering importance [TinyKittens 2015] or just to increase changes of finding a new home for homeless cats and kittens [TheCritter-Room 2015].

People are also more interested in health and safety of their own pets. Lower prices of video technologies and smartphone popularity helped to create various apps for controlling your pet lifestyle. [ZilliansInc 2015] created an automatic feeder with cat identification system and connected it to the smartphone app (product can be seen in fig. 1. It allows owner to supervise the amount of dry food and water his cats are consuming everyday, to receive alarms in case of abnormalities and even spy on them while eating.

The PiP [ThePetRecognitionCompany 2015] app is designated to help people find their lost cats and dogs. According to American Human Society, almost 3.5 million pets are lost each year and about 80% of dogs and 98% of cats are never re-united with their families. The authors of the smartphone app allows the owners to made profiles for their pets (including pictures) and in case the pet is missing, you



Fig. 1: CatFi product used by a cat.

launch the "alarm" on that pet, together with last seen location. If anyone finds a homeless animal, he may send the photo of it to the PiP app, which automatically try to match one of the lost pets with found pet.

Increasing number of pets' photographies also encouraged camera producers to implement appropriate support for it in their devices. Right now, Fujifilm company introduced Auto Dog / Cat Detection function [FujiFilm 2015] in some of their cameras that allows to automatically detect pet's face on the image and put auto-focus on it.

1.1 Problem description

Almost all of the applications and initiatives described in previous section already use or would benefit from the cats detection and/or identification. In this research I would like to focus on topic of cat detection in images. The main interest is on the use case of detecting cats presence in shots from live cameras and streams. Such images are characterized by diversified, often bad lighting settings, with possible presence of another animals or human beings. I will perform a state-of-the-art survey and then perform a comparison of existing, available algorithms.

2. THE STATE OF THE ART

Reliable comparison of existing software for detecting cats in pictures is hard to performed. It is caused by a few factors:

- **black-box algorithms** - many algorithms, that claim to have high accuracy, are part of the hardware product [ZilliansInc 2015] or closed-source app [ThePetRecognitionCompany 2015]. This way it is impossible to make a theoretical comparison of the algorithms and usually significantly hampers possibility of performing experiments with test data, as the products' interfaces are not adjusted to bulk operations.
- **high diversity of use cases** - although there are a lot of studies focusing on cat face detection, usually the classification or testing environmental does not check robustness of the algorithm in typical real-time scenarios. The classification is e.g. made as cat vs. still life [Fleuret and Geman 2007; Akihiko Yamada et al. 2011] or as cat/dog breed detection [Omkar M. Parkhi et al. 2012].
- **different test sets** - even for algorithms with similar use cases, e.g. cat vs. still life, different test set are in use to estimate the accuracy of algorithm, which makes it significantly harder to perform direct comparisons. Usually articles use their own data sets [Omkar M. Parkhi et al. 2012; Fleuret and Geman 2007; Akihiko Yamada et al. 2011] and only some of them compare to standard data bases, like Pascal 2007 cats database [Omkar M. Parkhi et al. 2012].

As the main idea of this research is to find an algorithm that will allow for highly accurate detection of cat faces in the pictures, only articles presenting "cat vs. still life" classification will be taken under consideration.

2.1 Unsupervised image categorization

One of the more interesting approaches to this problem was presented in [Le et al. 2012]. Researchers from Google and Stanford university constructed large scale unsupervised learning tool. It was a 9-layered neural network build on a cluster

that consisted of 1 000 machines with total of 16 000 cores. The network was trained with 10 million unlabelled images, all of them had the same size and local contrast normalization was applied on them. The model consisted of 1 billion connections and it took 3 days to train the classifiers. Final results are very promising for overall image categorization issue, because authors claims to have 15.8% accuracy for ImageNet dataset (22 000 object categories), which is a 70% improvement in comparison to the previous state of the art.

However, the most interesting part of [Le et al. 2012] for this research is that after performing face-detection experiment, authors tried to used trained classifiers to detect cats. The results are very promising - the random guess had 64.8% chance, while this algorithm was able to detect cats with 74.8% accuracy.

This algorithm and use of deep learning networks seems to be a promising technique for classification of images categories, however seemed to be lacking needed accuracy for binary classification (cat vs. not-cat).

2.2 Pet detection method for digital camera

Another interesting project was performed by [Akihiko Yamada et al. 2011]. They tried to construct a cat and dogs face detection mechanism for use in camera auto-focus and auto-exposure. They were not satisfied with current method, because camera has a very high constraints about calculation time and power consumption, were most of the studies focus on high recall and small false-detection. Authors decided to use four direction features as an input for the classifier. In order to speed-up the process, classification of an image window was hierarchical - if coarse classifier did not detect a sign of face in the window, it was automatically rejected and not send to more detailed and feature-rich fine classifier.

The results are very promising. As described in fig. 2, with the average recall on 90%, the false positive count is around 20. Although the FPC may seem high (the sample size was 400, with 200 of non-animal pictures), it is still on user-acceptable level. What is

more important, the processing time was only 245ms on ARM9 device with 166MHz core. It is impressive result, especially because this solution turned out to be also very robust to illumination changes. Authors measured how color and intensity fluctuation caused by illumination changes influence the recall rate. Intensity changes didn't cause the recall rate to drop under 80% even on very high intensity values. With 50% color fluctuation recall rate still oscillate between 80%-90%.

Although the results are not perfect (usually recall rate do not exceed 90%), short processing time allows to use this algorithm in real-time cases, which may be important case with live streaming observation. This solution can be used as a part of pet-identification algorithm, that informs appropriate person when the cat

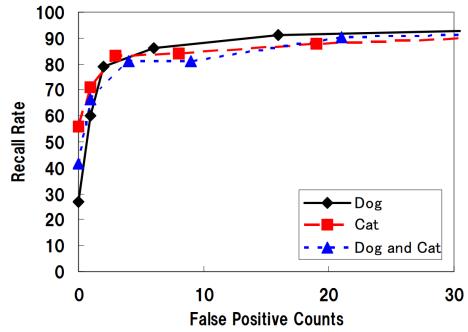


Fig. 2: ROC curves achieved as a result of [Akihiko Yamada et al. 2011].

appears in the camera (e.g. if it is known it is pregnant and needs to be caught and checked) and in such cases time is always a crucial value. It could be an interesting experiment to see if a little bit of power efficiency (i.e. better processing unit) could significantly improve the recall rate.

2.3 Shape and texture features for cat face detection

The solution presented in [Zhang et al. 2008] is the only one with existing available implementation. Also, as part of it, the on-line interface was provided for testing [Heather 2015]. Authors of [Zhang et al. 2008] decided, that both texture and shape are important features of cat face detection. Their solution has two steps - in first, they create separately trained shape and texture detector, both based on Haar-like features. As a second step they are jointly training fusion classifier (joint shape and texture). Both shape and texture classifier were trained on differently normalized data. Their features are based on oriented gradients.

Important part of that study was creating a database of 10 000 annotated cat pictures [Zhang et al. 2015] (annotations of 9 points in the cat face, like ears triangles, nose, mouth). This pictures will be the part of this research dataset.

The results of that study are presented in fig. 3. The recall rate oscillates around 90% with 100 FPC for the joint texture and shape detector. Because huge focus of [Zhang et al. 2008] was on feature engineering and comparing their features with different standard approaches, authors also made a direct comparison, using Pascal 2007 data set.

They claim that Average Precision (AP) of their approach is 0.364, when best reported method has AP equalled to 0.24.

The study seemed to find interesting and quite straight forward method of detecting cat's face. Coming from human face detection, they adjusted their feature choice with respect to differences appearing in cat face and then tested different type of gradient features to pick the most suitable one for this aim. However, it should be noted that this article is 7 years old and probably some of the studies presented in previous sections may have better results, but it is hard to estimate, as everyone are using different databases for testing.

2.4 Stationary features and cat detection

In the study [Fleuret and Geman 2007] authors presents a novel approach based on pose-indexed features. Firstly, the image is processed using custom edge detectors and grey scale histogram. Then, the second level feature is the head-belly pose, described as a reference frame, rotated accordingly to the pose vector. The expected result is a frame containing the cat.

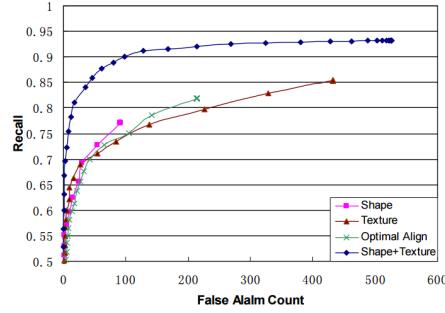


Fig. 3: ROC curves achieved as a result of [Zhang et al. 2008].

The classifiers used in that project are boosted trees and asymmetric weighting by sampling. The classifiers are using hierarchical search strategy, to avoid very costly scene processing. By using this methods, authors managed to avoid fragmentation of the data for separate training of pose classifiers.

The results of that study for the head-belly detection task are presented in fig. 4. The recall is around 80% and 85% (depending on method), with FAC equalled to respectively 17 and 11 on average with variation 3.5. Authors also performed head detection experiment with a little bit better results (recall 95% with FAC ranging from 6.65 to 12.31 depending on method).

Summing up, that novel approach seems to be very effective, but the images were mostly focusing on cat and still life cases, so it is hard to estimate the efficiency of this study in e.g. forest live streaming. Also, that study, similar to previous one described, is quite old if we take under consideration how fast machine learning is developing in last years, so it would probably be an interesting project to revise those ideas and try to apply them to modern approaches.

2.5 Detecting and tracking lions in videos

Very promising results and approaches are presented in study [Burghardt et al. 2006]. Authors are trying to create a tool for detecting the lion face in video footage, track it and based on the behaviour, add simple description (e.g. *"lioness walking towards camera, lies down on ground"*).

They are using Haag-like features, first presented in [Viola and Jones 2001]. During the studies, they tested that broader Haag-like feature set, like the one introduced in [Lienhart and Maydt 2002], may strongly increase accuracy of the algorithm. Then AdaBoost was used for finding the best classifier.

The result is tree classifier, trained to recognize frontal lion face. Then authors are using additional techniques to track the lion face movement, even when the pose is different (not frontal), but this is unfortunately out of scope of this project. However, what is interesting, is the final ROC curve, presented in fig. 5. The recall is equalled to 93% with false acceptance rate at the level of $10^{-4}\%$.

That study is offering effective solution to the problem described in this research. With satisfying level of accuracy and implementation based on standard solution ([Viola and Jones 2001]), even though the article is fairly old (9 years), it is probably the most promising technique.

2.6 Discussion

Described articles are presenting different approaches to the same problem, with focus on different factors (accuracy, efficiency, performance). In regards to the main

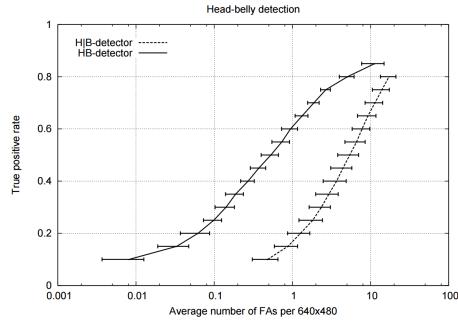


Fig. 4: ROC curves achieved as a result of [Fleuret and Geman 2007].

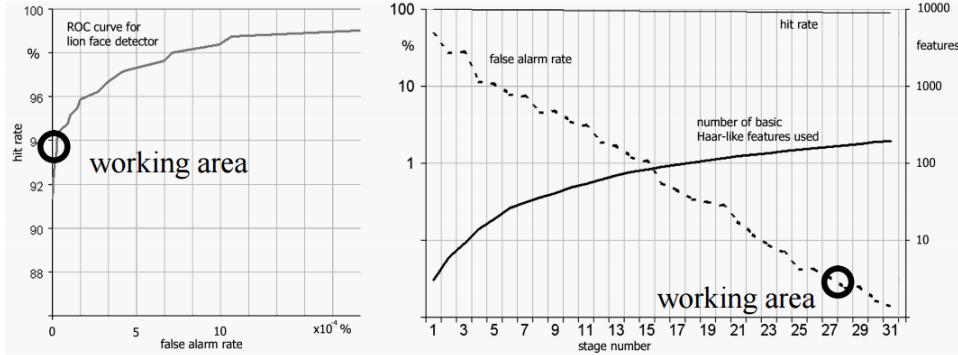


Fig. 5: (left) ROC curve for lion face training using 680 positive and 1000 negative images; (right) Outline of the learning process focused on the decrease of false alarms compared to the number of necessary single Haar-like features used to achieve this reduction, as described in [Burghardt et al. 2006].

problem in this article (cat detection in live streams), probably the studies described in section 2.2 and 2.5 are the best approaches. Those are the only studies that perform actual experiments focused on illumination invariance and/or live video samples. It is hard to estimate robustness of algorithm depending on pose variation, because none of the study actually focus on that, usually omitting this problem in some way (i.e. tracking techniques in [Burghardt et al. 2006]). Although everyone are mentioning that cat detection is harder than human detection because of much greater pose possibilities, not many tests are performed to check the robustness of algorithm to it, usually assuming that the classifier is trained only for the frontal view.

The main difficulty in comparing various methods of cat face detection is using different databases with different sizes (ranging from 400 to a million) and different non-cat data. From all the studies, only the one described in sec. 2.1 did have a database containing human faces or other animals. Interesting test was performed by an artists(!) group Shinseungback Kimyonghun [ShinseungbackKimyonghun 2015]. They used Kittydar [Heather 2015], the application based on algorithm described in sec. 2.3 and OpenCV face detection algorithm to parse pictures of humans and cats. They

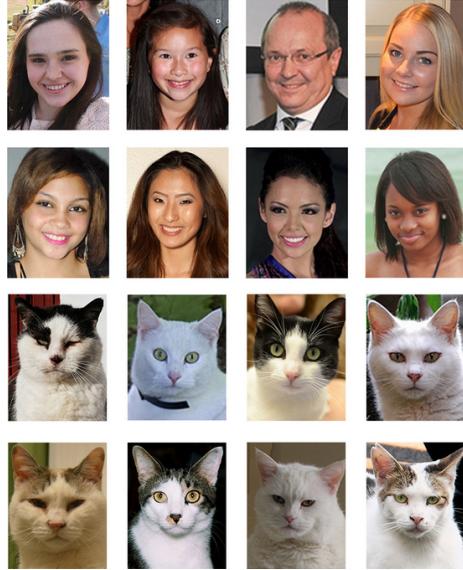


Fig. 6: First two rows - human recognized as cats by [Heather 2015]. Last two rows - cats recognized as human by OpenCV [Itseez 2015] face detection algorithm

ended up with 200 images, half of them with humans recognized by Kittydar as cats and half of them with cats recognized by OpenCV as humans. Part of their results is shown in figure 6.

3. PROBLEM ANALYSIS

This section is meant to describe the properties of algorithm suitable for cat detection. It will be started with recalling of expected algorithm applications:

- **be first part of cat identification system** - as such, it should have low false rejection rate. High false acceptance is allowed, as it would be filtered by identification subsystem.
- **work in poor light settings** - The live stream takes place in different locations - house, forests, gardens, often 24/7, and there may be situation when no artificial light is provided. The algorithm should be as robust as possible to light variance.
- **distinguish between cats and cat-like, humans, other animals** - a lot of current studies does not take under consideration that human faces or other animals may be the part of background. It is important to prepare features so that not too many false acceptance are caused by another beings (although, as mentioned above, false acceptance ratio is not the most important statistic).

3.1 Algorithm properties

With these applications in mind, one can summarize it with list of algorithm properties:

- robust to light/pose variation
- as low as possible false rejection count
- false acceptance count may be high
- additional focus on distinguishing between cats and other living creatures

Following considerations from [Zhang et al. 2008], it is good idea to discuss different methods that can be used in the solution:

- **sliding window detection vs. part based detection** - there are two different ways of detecting objects in the images. Sliding window scan all possible sub-windows in the image and try to classify each of them (does it have a face in it?). Part-based detects simple part of the objects (e.g. for human face recognition it could be eyes and mouth) and try to assemble them, following geometric constraints. Both of those methods are used, with successes, in human face detections. Probably sliding window would be a better method for this use case, more robust to illumination changes.
- **image features vs. gradient features** - image features are for example intensity values or PCA/wavelet coefficients. On the other side, gradient features are extracted from oriented gradients or the edge map. Both of this approaches demonstrated satisfying results in previous studies. Although image features are said to be less robust to illumination changes, they're successfully used in detection from video [Burghardt et al. 2006], so it is hard to decide one way or another before further experiments.

4. EXPERIMENTS

The experiments will be performed using Python interface of OpenCV library.

OpenCV [Itseez 2015] is an open source computer vision framework, developed from 1999, right now maintained by Itseez company. The library is cross-platform and has C, C++, Java, Python and even Matlab/Octave interfaces. It is released under BSD licence, which means it is free for both academic and commercial use. OpenCV is written in C++ and put strong focus on computational efficiency, supporting multi-core and GPU processing. It offers broad range of applications from computer vision, like facial recognition system, motion tracking, 2D/3D feature toolkits and also some machine learning techniques: different type of classifiers (SVM, Bayes, Random Forest, ANN) or boosting methods (boosting, gradient boosting). It has large user community (around 50 thousand) and is considered to be the best free tool for computer vision.

Python is general-purpose, high-level programming language. It supports object oriented programming and has elements of imperative, functional and procedural style. Python has dynamic typing and garbage collector. There are many different implementations of Python. The most popular and reliable one are also free and open-source. It can be packaged to be used on almost all the platforms. Its characteristic is simplicity and high reliability. Python has a lot of frameworks and libraries that allow it to be useful in many different fields, e.g. Django (web development), SciPy, SciKit (scientific and machine learning), OpenStack (system administration). Also, many existing frameworks are ported or have interfaces in Python (like OpenCV).

4.1 Tested algorithms

Algorithms, that will be tested in the experiments, are classifiers based on Haar-like features. Two approaches will be tested:

- *predefined cascade classifier* - OpenCV offers two existing, pre-trained frontal cat face classifiers.
- *trained cascade classifier* - OpenCV offers possibility of self-training the classifier with the data.

4.2 Classifier training

Following properties of classifier (and its influence on accuracy) will be tested during the experiment:

- *feature size* (how many Haar features are taken under consideration) - tested values will be: basic (only upright features), core, all (upright + 45 degree rotated feature set)
- *maximal false acceptance rate* (per stage).
- *number of stages* (tested values: 1, 5, 10).

Feature type (Haar or LBP - local binary patterns [Liao et al. 2007]) won't be considered as a changing property, because LBP usually has worse results (although it is also faster to learn). Number of positive and negative samples is constrained by out of memory exception on values only around 300-500 images, depending on rest of the factors, so also cannot be accuracy influencing factor (as it is already small).



Fig. 7: Images used for testing the classifiers - from left to right Picture 1, 2 and 3.

Boosting classifier is Gentle AdaBoost (used in some of the articles in previous section).

4.3 Training data set

The positive data set contains 800 photos of kittens from database introduced in [Zhang et al. 2008]. The negative data set has 800 photos of house interiors from SUN database [MIT 2015].

4.4 Testing

Because each of resulting classifier will have predefined accuracy, the method of validation will be based on testing each classifier with three images (fig. 7), that do not appear in training set, and subjectively estimating the results with scale consisting of 6 values:

- 0 - no part of cat has been found
- 1 - some parts of cat have been found, less than 20% found areas belongs to cat
- 2 - some parts of cat have been found, between 20-50% found areas belongs to cat
- 3 - some parts of cat have been found, between 50-80% found areas belongs to cat
- 4 - more than 80% found areas belongs to cat
- 5 - all found areas belongs to cat

4.5 Experiment - predefined cascades

As a first experiment, two predefined cascade classifiers from OpenCV were taken - *haarcascade_frontalcatface* and *haarcascade_frontalcatface_extended*. The difference between the is that standard one is using only basic set of Haar features, while extended use all of them.

Results of experiment are presented in tab.I. Although average accuracy of standard classifier seems to be higher than extended, it is caused by number of detected areas. Extended classifier apparently has higher activation threshold, which is causing rather not detecting the cat than making a mistake, but all the detected areas are always on the cat body/face.

Cascade	Picture 1	Picture 2	Picture 3	Average no. of areas
Standard	3	5	4	16-21
Extended	0	0	5	1

Table I: Results of testing the quality of predefined cascades

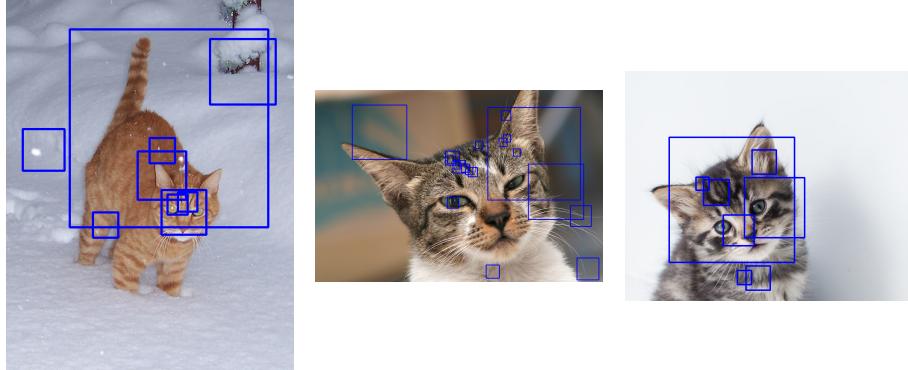


Fig. 8: Results of testing classifier 0.2B10.

4.6 Experiment - trained cascade

The second experiment will show difference and possibilities of simply classifier' training. The images used for testing the classifiers were only changed to the grey scale and hasn't been processed expect for that. The training sample size are:

- 300 positive, 500 negative for basic features
- 250 positive, 350 negative for core features
- 150 positive, 250 negative for all features

The recall is always around 99%.

The results are presented in tab. II. It can be observed, that the highest average evaluations are for classifiers with FAR around 0.1-0.2, with 10 stages and basic or core set of features (e.g. fig. 8). When classifier is searching for all features, a lot more areas are found, that are usually not even part of carts.

5. DISCUSSION

Although the classifiers were not able to detect cat face/body as one rectangle, smoothing the result should result in satisfying output. Probably preparing dedicated machine for training and increasing training set sizes (especially negative training set) would make algorithm more robust to different backgrounds.

6. CONCLUSIONS

Experiments have given interesting, but also expected results. Too high or too low false acceptance rate is causing over or underfitting. The number of stages is important part of algorithm and should be around 10. The feature set is important attribute, but not all of the Haar features are required for accurate cat detection. The classifiers can be trained easily to detect cats with satisfying precision. Previously describe articles (especially [Burghardt et al. 2006] and this experiments

Cascade	Picture 1	Picture 2	Picture 3	Average	Average no. of areas
0.1B1	1	2	1	1.5	25-75
0.1C1	3	4	1	2.3	12-70
0.1A1	1	2	3	2	4-13
0.1B5	4	4	1	3	16-45
0.1C5	3	3	4	3.3	30-75
0.1A5	2	3	1	2	113-130
0.1B10	0	5	5	3.3	0-2
0.1C10	2	5	5	4	2-8
0.1A10	5	3	5	4.3	2-25
0.2B1	3	3	1	2.3	10-59
0.2C1	3	3	1	2.3	8-68
0.2A1	2	5	3	3.3	1-11
0.2B5	4	4	2	3.3	47-89
0.2C5	3	4	3	3.3	57-93
0.2A5	1	2	2	1.6	97-148
0.2B10	4	4	5	4.3	8-19
0.2C10	4	4	5	4.3	5-21
0.2A10	2	3	5	3.3	5-25
0.4B1	4	2	1	2.3	7-55
0.4C1	3	3	1	2.3	7-50
0.4A1	2	3	4	3	3-7
0.4B5	2	4	3	3	84-122
0.4C5	2	3	2	2.3	115-129
0.4A5	2	3	3	2.6	87-134
0.4B10	3	4	2	3	28-54
0.4C10	4	4	2	3.3	28-78
0.4A10	2	3	2	2.3	75-123

Table II: Results from self-trained classifiers. In the classifier name, the first number means maximal false acceptance rate at each stage, the letter - feature set (basic, core, all) and the last number - number of stages.

are showing, that it is possible to train classifier based on Haar-like features to be a valuable tool for detecting cats on photos with different light setting and even slightly robust to pose.

7. BIBLIOGRAPHY

REFERENCES

- AKIHIKO YAMADA, KAZUHIRO KOJIMA, J. K., OKAMOTO, M., AND MURATA, H. 2011. Directional edge-based dog and cat face detection method for digital camera. *IEEE International Conference on Consumer Electronics*.
- BURGHARDT, T., ACL, J., AND THOMAS, B. T. 2006. Tracking animals in wildlife videos using face detection.
- FLEURET, F. AND GEMAN, D. 2007. Stationary features and cat detection. *Idiap Research Report*.
- FUJIFILM. 2009 (accessed June 22, 2015). *Auto Dogs & Cats Detection*. http://www.fujifilm.com/products/digital_cameras/pet/.
- HEATHER, A. 2013 (accessed June 22, 2015). *Kittydar - face detection for cat*. <http://harthur.github.io/kittydar/>.
- ITSEEZ. 2015 (accessed June 22, 2015). *OpenCV*. <http://opencv.org/>.
- ANGLEYANIMALPROTECTIONSOCIETY. 2015 (accessed June 22, 2015). *Trap-Neuter-Return (TNR) Pilot*. <http://www.tinykittens.com/cats>.

- LE, Q. V., RANZATO, M. A., MONGA, R., DEVIN, M., CHEN, K., CORRADO, G. S., DEAN, J., AND NG, A. Y. 2012. Building high-level features using large scale unsupervised learning.
- LIAO, S., ZHU, X., LEI, Z., ZHANG, L., AND LI, S. Z. 2007. Learning multi-scale block local binary patterns for face recognition.
- LIENHART, R. AND MAYDT, J. 2002. An extended set of haar-like features for rapid object detection. MIT. 2015 (accessed June 22, 2015). *Sun Database*. <http://groups.csail.mit.edu/vision/SUN/>.
- OMKAR M. PARKHI, A. V., ZISSERMAN, A., AND JAWAHAR, C. V. 2012. Cats and dogs.
- SHINSEUNGBACKKIMYONGHUN. 2013 (accessed June 22, 2015). *Cat or Human*. http://ssbkyh.com/works/cat_human/.
- THECRITTERROOM. 2015 (accessed June 22, 2015). *Foster Kitten Cam*. <https://livestream.com/FosterKittenCam>.
- THEPETRECOGNITIONCOMPANY. 2013 (accessed June 22, 2015). *PiP - reuniting lost pets with their families*. <http://www.petrecognition.com/>.
- TINYKITTENS. 2015 (accessed June 22, 2015). *Rescue Kitten Reality TV*. <http://www.tinykittens.com/>.
- VIOLA, P. AND JONES, M. 2001. Robust real-time object detection.
- ZHANG, W., SUN, J., AND TANG, X. 2008. Cat head detection - how to effectively exploit shape and texture features. *Proc. of European Conf. Computer Vision* 4, 802–816.
- ZHANG, W., SUN, J., AND TANG, X. 2008 (accessed June 22, 2015). *Cat database*. <http://137.189.35.203/WebUI/CatDatabase/catData.html>.
- ZILLIANSINC. 2015 (accessed June 22, 2015). *CatFi - Smart Feeder*. <http://catfi.com/>.

Received June 2015;