

Impact Factor:

ISRA (India) = 3.117
ISI (Dubai, UAE) = 0.829
GIF (Australia) = 0.564
JIF = 1.500

SIS (USA) = 0.912
PIHHI (Russia) = 0.156
ESJI (KZ) = 8.716
SJIF (Morocco) = 5.667

ICV (Poland) = 6.630
PIF (India) = 1.940
IBI (India) = 4.260
OAJI (USA) = 0.350

SOI: [1.1/TAS](#) DOI: [10.15863/TAS](#)

International Scientific Journal Theoretical & Applied Science

p-ISSN: 2308-4944 (print) e-ISSN: 2409-0085 (online)

Year: 2019 Issue: 06 Volume: 74

Published: 15.06.2019 <http://T-Science.org>

QR – Issue



QR – Article



Konstantin Sergeevich Chebykin

Peter the Great St.Petersburg Polytechnic University
student

chebykinkostya@gmail.com

Vadim Andreevich Kozhevnikov

Peter the Great St.Petersburg Polytechnic University
Senior Lecturer

vadim.kozhevnikov@gmail.com

RECOGNITION OF BODY PART USING NEURAL NETWORKS

Abstract: This work belongs to the training of neural networks for solving problems related with computer vision. The first chapter discusses the general concepts of neural networks and their structure. The second chapter reviews the models of convolutional neural networks and the ways of learning them. The third chapter describes the project architecture, implementation and testing.

Key words: machine learning, computer vision, convolutional neural networks, recognition of images

Language: English

Citation: Chebykin, K. S., & Kozhevnikov, V. A. (2019). Recognition of body part using neural networks. *ISJ Theoretical & Applied Science*, 06 (74), 233-239.

Soi: <http://s-o-i.org/1.1/TAS-06-74-27> **Doi:**  <https://dx.doi.org/10.15863/TAS.2019.06.74.27>

Introduction

Recently, tasks related to computer vision are very popular. With the help of machine learning you can detect, explore, recognize any object or help the robot to navigate in the area. Every year the number of tasks in this sphere grows and new opportunities in the field of computer vision are opening up. Pattern recognition is a scientific discipline whose goal is to identify objects according to several criteria or classes [1].

Where could this come in handy? Now there are several interesting examples of these solutions. On December 5, 2016, Amazon opens its first store without AmazonGo sellers and by 2019 opens 4 such stores [2]. Shops are operating with the following technologies: computer vision, various sensors and deep learning. Also more uses can be found in video games of the virtual reality. Players will not have to keep gamepads constantly in their hands, it will be enough to use only hands to play without gamepads.

This solution also can be used in unmanned vehicles. Every unmanned vehicle has a video camera that helps to identify obstacles such as pedestrians or cyclists. People can warn the car with any gestures. This can save the life of a passenger or other people.

Accordingly, the goal is to study and learn neural networks and their application for solving problems of searching for human joints. To achieve this goal, we formulate the following tasks [3]:

- 1) Study the problem of pattern recognition;
- 2) Study of the features of the applications of neural networks;
- 3) Search and data processing for training;
- 4) Training and optimization of models;
- 5) Evaluation of the results.

The aim of the article

The aim of the current work is the implementation of neural networks that can predict the coordinates of human joints. Neural networks can be used in unmanned vehicles to provide protection to passengers and pedestrians. And before this aim, the following tasks were set: analyze the use of neural networks in the field of computer vision, collect the necessary data for training, implement and train neural networks using the Python language, implement a simple application to visualize the results.

Stack of used technologies

Impact Factor:

ISRA (India)	= 3.117	SIS (USA)	= 0.912	ICV (Poland)	= 6.630
ISI (Dubai, UAE)	= 0.829	PIHHI (Russia)	= 0.156	PIF (India)	= 1.940
GIF (Australia)	= 0.564	ESJI (KZ)	= 8.716	IBI (India)	= 4.260
JIF	= 1.500	SJIF (Morocco)	= 5.667	OAJI (USA)	= 0.350

Today there are many different languages for working with neural networks, such as: Python, C++, R, Java and others, as well as various libraries and frameworks, for example, Coffee, Keras, Tensorflow, Theano, Caffe.

The first advantage of Python and R is an uncomplicated syntax, it is “elegant” on the one hand, and “mathematical” on the other, its semantics have a special correspondence to many common mathematical ideas. And languages like C++ and Java are designed for maximum performance, but with a complex syntax for understanding and writing.

The second advantage of Python is the prevalence among developers. Since it is the most common language, this implies the availability of many different libraries and machine learning tools, which are mostly oriented towards working with Python. Tensorflow, Keras - one of the most popular libraries and frameworks, numpy, matplotlib, jupyter notebook - one of the most popular tools for work.

Tensorflow is the most common open source library for numerical calculations using data flow graphs [4]. This library is cross-platform, so it can work on a GPU, CPU, and also on TPU tensor processors. But one of the main advantages is that there are many implemented architectures for various neural networks.

Keras is a framework written in Python that can work on top of Tensorflow, Teano, and other libraries for working with neural networks. It was implemented to improve interaction speed with neural networks. The ability to move from idea to result with the least possible delay is the key to conducting good research [5].

As we are trying to achieve maximum precision and convenience, Python has been chosen. Tensorflow, Keras, numpy, matplotlib, jupyter notebook are also used in the solution.

To implement the client-user application, the Qt graphical framework for the Python programming language was chosen.

Project Architecture

To implement the project, it was decided to use two consecutive neural networks. The task of the first neural network is to specify the bounding box for the people in the image. The task of the second neural network is to predict the coordinates of the joints for people in the bounding boxes.

Model overview for the first convolutional neural network

The task of the first neural network is to classify people in the image. To solve this problem, there are many models, so we consider them and choose the best in our opinion. We will consider the model in chronological order.

Regions Convolutional Neural Networks (R-CNN)

In 2014, a small team at the University of California at Berkeley publishes a neural network that can detect objects in an image [6]. Object detection is the task of searching for various objects and their classification. For example, to recognize a cat or a dog in an image. The purpose of R-CNN is to take images and correctly predict bounding box of the objects. A bounding box is a rectangle that bounds the shape of a more complex geometric model and is determined by the coordinates of the upper left and right lower points.

1. An arbitrary resolution image is supplied to the R-CNN input. The original image is divided into rectangles of different sizes with the help of SelectiveSearch, and they are called object candidates. 2000 candidates are predicted for each image.

2. Each candidate is given for a resolution of 224x224.

3. Next, using their own implementation or the implementation of Krizhevsky CNN, determine the weights of features of the 4096-dimensional vector for each candidate.

4. Using the classifier (SVM) recognize the object by signs.

5. Using the regressor, we predict bounding box of the objects.

SelectiveSearch - an algorithm for determining the similar regions. It is a hierarchical grouping of similar areas based on compatibility of shape, color, structure, and size. In other words, this algorithm is for clustering intersected regions in an image [7].

Fast R-CNN

The usual R-CNN had 2 big problems. 2000 candidates were predicted for each image and each candidate was processed by a neural network. Accordingly, the first problem was the speed of learning and testing. The second is that it was necessary to train CNN, the classifier (SVM), which recognizes the object, and the regression model, in order to narrow the bounding box for the image.

Problem solving was proposed in 2015 by Ross Hirschik, who participated in the development of R-CNN, and named the model Fast R-CNN [8].

The author solves the first problem. He is not applying CNN for each candidate, just apply only to the original image, and then use the RoI Pooling layer to see if one of the 2000 areas is suitable for the object. Pooling receives data in the form of a C x H x W map from the last CNN convolutional layer, as well as the height(h) and the width(w) of the challenger. They are compared, and a decision is made whether the applicant is suitable for the object — if so, it is inserted into the layer of the ROI pool.

Impact Factor:

ISRA (India) = 3.117	SIS (USA) = 0.912	ICV (Poland) = 6.630
ISI (Dubai, UAE) = 0.829	PIHHI (Russia) = 0.156	PIF (India) = 1.940
GIF (Australia) = 0.564	ESJI (KZ) = 8.716	IBI (India) = 4.260
JIF = 1.500	SJIF (Morocco) = 5.667	OAJI (USA) = 0.350

The second problem was solved by the fact that the regressor, the classifier and CNN were combined into one model and trained together.

Faster R-CNN

In spite of all the advantages of Fast R-CNN, there were still some minor problems, and they were connected with the candidates. In 2016, under the leadership of Juan Sun, the main researcher at Microsoft Research, and his team, a new model Faster R-CNN was proposed. The idea of the new model was to use the predicted areas for several candidates at once, instead of one candidate. In this model, this approach was implemented and it turned out to be quite successful, which made it possible to increase accuracy.

YOLO

YOLO or You only look once is a model for detecting an object, which is very different from the previously considered models. In YOLO, one convolutional network predicts bounding boxes and class probabilities for each box. How YOLO works - we take an image and divide it into an SxS grid, we divide each cell block into M parts. For each part, the neural network predicts the probability class and the offset values for the bounding box. A box having a class probability above the threshold value is selected and used to determine the position of the object in the image. The class probability is calculated as follows (1):

$$P = P(obj) \cdot IOU \quad (1)$$

where $P(obj)$ - the probability that the bounding box contains an object, IOU - the ratio of the region of overlap of the predicted and true bounding rectangles to the region of the union of these rectangles

SSD

Further development of the YOLO idea was developed in the SSD model, that using the same principle of object detection all over image that is reflected in its name - Single Shot Detector. The SSD model is the first deep network for detecting objects that does not use signs consisting of bounding boxes.

It leads to a significant increase in performance and detection of objects with higher accuracy. In this model, they abandoned the use of bounding boxes, which resulted in an increase in the speed of the SSD. Improving accuracy occurs by using a small convolutional filter to predict the class of an object and a convolutional filter to correct the position of bounding boxes with different aspect ratios and sizes. These convolutional filters are applied to several feature maps, both in the early and in the later layers of the network, which makes it possible to detect objects of various sizes.

YOLOv2

YOLOv2 is the second version of YOLO, the purpose of which is to significantly improve accuracy while accelerating. One of the main differences is that a normalization layer is added to each convolutional layer. This method lies in the fact that some layers of the neural network are fed to the input data, pre-processed and having zero expectation and unit variance. Another major difference is that we replace all fully connected layers with convolution layers, which allows us to process large images and facilitates the training of the classifier.

Model YOLOv2 with a resolution of 416x416 was chosen to implement the first neural network. Since the ratio of accuracy to learning speed of this model is maximum. This is demonstrated in the article "YOLO9000: Better, Faster, Stronger", which compared the accuracy and speed of training models that were trained on PASCAL VOC 2007 data [10]. Consider the HC1 architecture in more detail (Figure 1). The neural network consists of 6 convolutional, 6 pooling layers, 8 batch normalization layers and 3 convolutional layers at the end. All convolutional layers, except the last, have the leaky ReLU activation function, and the last layer has a linear activation function. Pooling layers are max-pooling layers with the choice of the maximum element.

Source	Train?	Layer description	Output size
		input	(?, 416, 416, 3)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 416, 416, 16)
Load	Yep!	maxp 2x2p0_2	(?, 208, 208, 16)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 208, 208, 32)
Load	Yep!	maxp 2x2p0_2	(?, 104, 104, 32)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 104, 104, 64)
Load	Yep!	maxp 2x2p0_2	(?, 52, 52, 64)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 52, 52, 128)
Load	Yep!	maxp 2x2p0_2	(?, 26, 26, 128)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 26, 26, 256)
Load	Yep!	maxp 2x2p0_2	(?, 13, 13, 256)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 13, 13, 512)
Load	Yep!	maxp 2x2p0_1	(?, 13, 13, 512)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 13, 13, 1024)
Init	Yep!	conv 3x3p1_1 +bnorm leaky	(?, 13, 13, 1024)
Init	Yep!	conv 1x1p0_1 linear	(?, 13, 13, 30)

Fig. 1 First neural network architecture

Impact Factor:

ISRA (India)	= 3.117	SIS (USA)	= 0.912	ICV (Poland)	= 6.630
ISI (Dubai, UAE)	= 0.829	PIHHI (Russia)	= 0.156	PIF (India)	= 1.940
GIF (Australia)	= 0.564	ESJI (KZ)	= 8.716	IBI (India)	= 4.260
JIF	= 1.500	SJIF (Morocco)	= 5.667	OAJI (USA)	= 0.350

Model overview for the second convolutional neural network

The task of the second neural network is to predict the coordinates of the joints based on the results of the first neural network. Changing the number of different layers and functions of activations, the model is presented below (Figure 2). At the entrance, the neural network receives an image with a resolution of 224x224. Our model consists of: 5 convolutional layers with 2 convolution cores 11x11 and 3 3x3 cores, 3 max-pooling layers with a 2x2 pool

core, 4 fully connected layers with the number of neurons 4096, 4096, 1000, 28. Each convolutional and max-pooling layers have activation function ReLU, batch normalization layer. Each fully connected layer contains the activation function tanh, batch normalization layer and dropout layer with a value of 0.3 so that the model is not overfitting.

The last layer contains 28 neurons, each of which contains the position of x or y relative to the center. The value of the neuron in the output layer takes the value [-0.5;0.5].

Layer (type)	Output Shape	Param #
conv2d_1 (Conv2D)	(None, 54, 54, 96)	34944
activation_1 (Activation)	(None, 54, 54, 96)	0
max_pooling2d_1 (MaxPooling2D)	(None, 27, 27, 96)	0
batch_normalization_1 (Batch Normalization)	(None, 27, 27, 96)	384
conv2d_2 (Conv2D)	(None, 17, 17, 256)	2973952
activation_2 (Activation)	(None, 17, 17, 256)	0
max_pooling2d_2 (MaxPooling2D)	(None, 8, 8, 256)	0
batch_normalization_2 (Batch Normalization)	(None, 8, 8, 256)	1024
conv2d_3 (Conv2D)	(None, 6, 6, 384)	885120
activation_3 (Activation)	(None, 6, 6, 384)	0
batch_normalization_3 (Batch Normalization)	(None, 6, 6, 384)	1536
conv2d_4 (Conv2D)	(None, 4, 4, 384)	1327488
activation_4 (Activation)	(None, 4, 4, 384)	0
batch_normalization_4 (Batch Normalization)	(None, 4, 4, 384)	1536
conv2d_5 (Conv2D)	(None, 2, 2, 256)	884992
activation_5 (Activation)	(None, 2, 2, 256)	0
max_pooling2d_3 (MaxPooling2D)	(None, 1, 1, 256)	0
batch_normalization_5 (Batch Normalization)	(None, 1, 1, 256)	1024
flatten_1 (Flatten)	(None, 256)	0
dense_1 (Dense)	(None, 4096)	1052672
activation_6 (Activation)	(None, 4096)	0
dropout_1 (Dropout)	(None, 4096)	0
batch_normalization_6 (Batch Normalization)	(None, 4096)	16384
dense_2 (Dense)	(None, 4096)	16781312
activation_7 (Activation)	(None, 4096)	0
dropout_2 (Dropout)	(None, 4096)	0
batch_normalization_7 (Batch Normalization)	(None, 4096)	16384
dense_3 (Dense)	(None, 1000)	4097000
activation_8 (Activation)	(None, 1000)	0
dropout_3 (Dropout)	(None, 1000)	0
batch_normalization_8 (Batch Normalization)	(None, 1000)	4000
dense_4 (Dense)	(None, 28)	28028
activation_9 (Activation)	(None, 28)	0
Total params: 28,107,780		
Trainable params: 28,086,644		
Non-trainable params: 21,136		

Fig. 2 Second neural network architecture

Impact Factor:

ISRA (India)	= 3.117	SIS (USA)	= 0.912	ICV (Poland)	= 6.630
ISI (Dubai, UAE)	= 0.829	PIHHI (Russia)	= 0.156	PIF (India)	= 1.940
GIF (Australia)	= 0.564	ESJI (KZ)	= 8.716	IBI (India)	= 4.260
JIF	= 1.500	SJIF (Morocco)	= 5.667	OAJI (USA)	= 0.350

Neural network training

The first neural network was trained on 2000 images containing less than 4 people. All bounding rectangles were selected manually. The number of learning epochs is set to 100. The number of epochs shows how many times a neural network will be trained on training data. As a function of the losses selected (2):

$$\sum_{i=1}^N \sum_{j=1}^M L_{i,j} \cdot ((x_i - x_j)^2 + (y_i - y_j)^2 + (h_i - h_j)^2 + (w_i - w_j)^2) \quad (2)$$

where $i=1..N$ - number of predicted boxes, $j=1..M$ - number of results boxes, x_i and y_i - coordinates of the upper left point of the bounding box, h_i - height of the bounding box, w_i - width of the bounding box, $L_{i,j}$ - confidence coefficient, which is determined by IoU, if

$\text{IoU} \geq 0.5$ then $L_{i,j} = 1$, else $L_{i,j} = 0$. IoU = intersection area divided by total area.

Adam was chosen as the optimizer for the first neural network, because we need an optimizer that is more adapted to a large amount of data with a lot of noise and more adapted to work with deep neural networks. The model was trained for 12 hours and after every 6 epochs it was tested on the tested data. The mAP metric was used as the accuracy metric. If the value of $\text{IoU} > 0.5$ for the predicted rectangle and the correct result, then we assume that the neural network predicted the correct result [9]. The graph of dependence of the mAP metric is presented below (Figure 3).

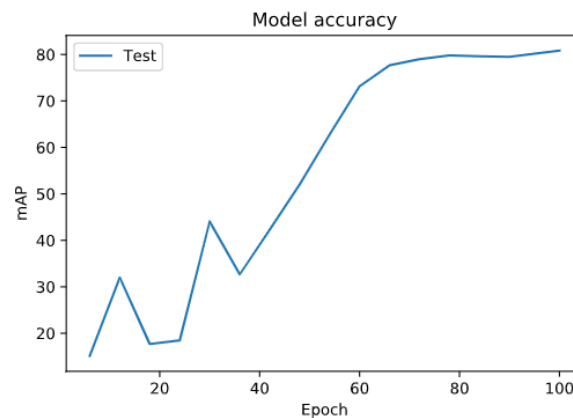


Fig. 3 Graph of the dependence of the mAP metric on the number of epochs

The result for the image that was not included in the training and the test data is presented below (Figure 4).

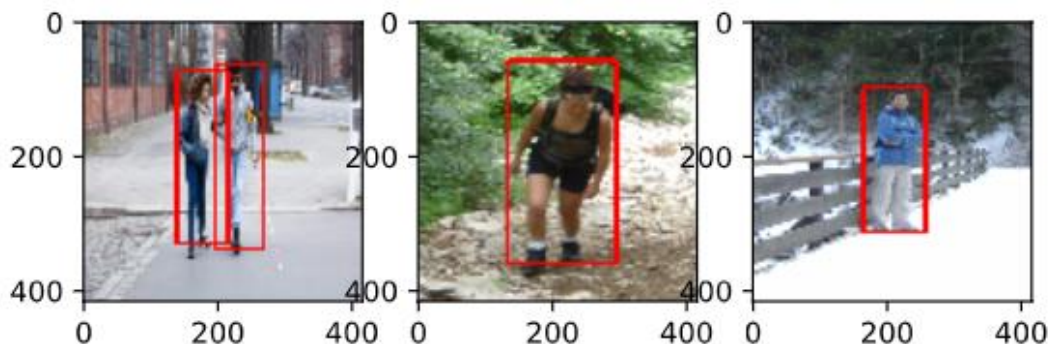


Fig. 4 The result of the first neural network

The second model was trained on LSP data (Leeds Sports Pose Dataset), which includes images of people involved in sports and the coordinates of the joints marked on them. The data were broken down in

relation to 80/20 for training and test data. The number of epochs is set to 50.

Adam was chosen as an optimizer. Since he showed himself well when learning the first neural

Impact Factor:

ISRA (India) = 3.117	SIS (USA) = 0.912	ICV (Poland) = 6.630
ISI (Dubai, UAE) = 0.829	PIHII (Russia) = 0.156	PIF (India) = 1.940
GIF (Australia) = 0.564	ESJI (KZ) = 8.716	IBI (India) = 4.260
JIF = 1.500	SJIF (Morocco) = 5.667	OAJI (USA) = 0.350

network. The MSE metric (3) was chosen as the estimation and loss function metric:

$$MSE = \frac{\sum_{i=1}^n (x_i - y_i)^2}{n} \quad (3)$$

where x_i - neural network result, y_i - true result.

The neural network has been trained for about 60 hours. During the training, it was possible to achieve 60% accuracy, which is explained by the fact that 28 parameters predict a neural network and rather complex training and training data (Figure 5).

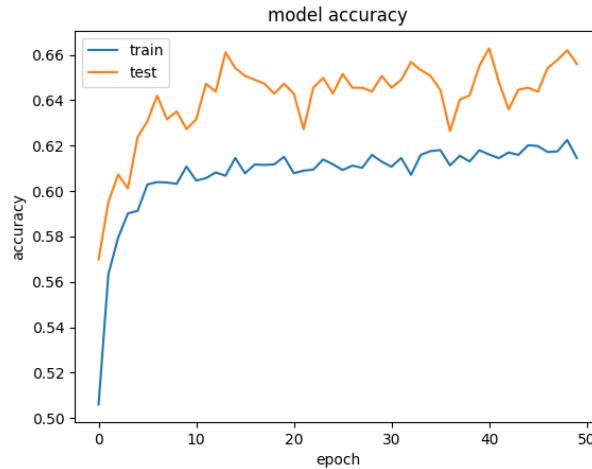


Fig. 5 Graph of the dependence of the accuracy metric on the number of epochs

For convenient use of neural networks, a window application has been written that has a convenient and simple interface (Figure 6). The windowing

application was written in python and using the Qt graphical framework.

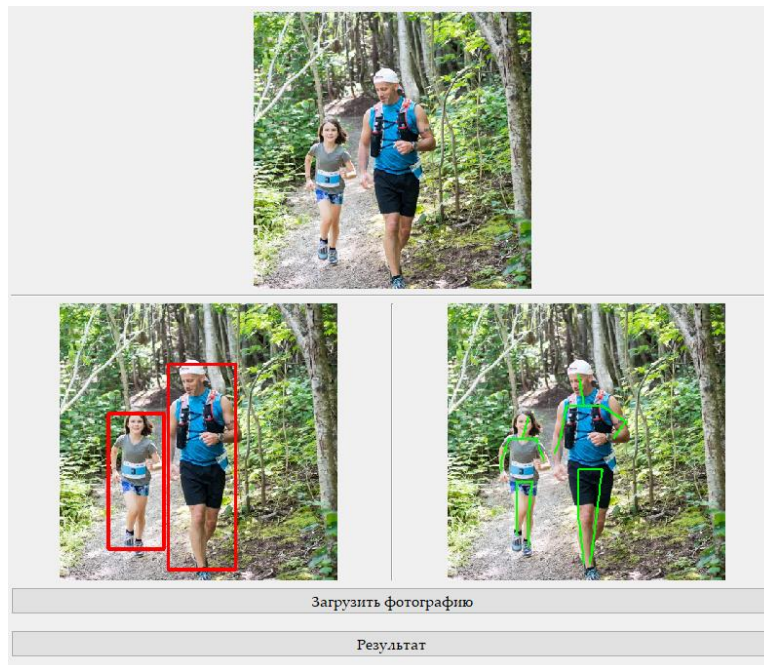


Fig. 6 Result of work window application

Conclusion

Two neural networks were trained, each responsible for its own task. The result of which is marked up people in the images, which can be useful in use in unmanned vehicles. As a continuation of the

work, it is possible to come up with a more complex architecture for the second neural network to increase accuracy. And in the future, to organize the prediction of the position of a person in 3D using several pictures from 2D.

Impact Factor:	ISRA (India) = 3.117	SIS (USA) = 0.912	ICV (Poland) = 6.630
	ISI (Dubai, UAE) = 0.829	PIHHI (Russia) = 0.156	PIF (India) = 1.940
	GIF (Australia) = 0.564	ESJI (KZ) = 8.716	IBI (India) = 4.260
	JIF = 1.500	SJIF (Morocco) = 5.667	OAJI (USA) = 0.350

References:

1. Chernogorov, Y. V. (2019). Pattern recognition methods. Retrieved April 04, 2019, from <https://moluch.ru/archive/132/36964>
2. (n.d.). The official website of the Amazon service [online]. Retrieved April 04, 2019, from <https://www.amazon.com/mestetskii04course.pdf>
3. (n.d.). Course of lectures. Mathematical methods of pattern recognition [online]. Retrieved April 04, 2019, from <http://www.ccas.ru/frc/papers/mestetskii04course.pdf>
4. (n.d.). What is the TensorFlow machine intelligence platform? [online] Retrieved April 04, 2019, from <https://opensource.com/article/17/11/intro-tensorflow>
5. (n.d.). The official documentation site of the Keras library [online]. Retrieved May 04, 2019, from <https://keras.io/>
6. (n.d.). Object detection: speed and accuracy comparison. [online]. Retrieved May 04, 2019, from https://medium.com/@jonathan_hui/object-detection-speed-and-accuracy-comparison-faster-r-cnn-r-fcn-ssd-and-yolo-5425656ae359
7. (n.d.). Selective Search for Object Recognition [online]. Retrieved May 04, 2019, from <http://www.huppelen.nl/publications/selectiveSearchDraft.pdf>
8. (n.d.). Fast R-CNN [online]. Retrieved May 04, 2019, from <https://arxiv.org/pdf/1504.08083.pdf>
9. (n.d.). Course of lectures. Selective Search for Object Recognition. [online]. Retrieved May 04, 2019, from <http://www.cs.cornell.edu/courses/cs7670/2014sp/slides/VisionSeminar14.pdf>
10. (n.d.). YOLO9000: Better, Faster, Stronger. [online]. Retrieved May 04, 2019, from <https://arxiv.org/pdf/1612.08242.pdf>