

Deep Reinforcement Learning

SoftServe Data Science Group

Ihor Kostiuk, Roman Grubnyk, Iurii Milovanov, Ievgenii Baliuk, Pavlo Kramarenko, Tetiana Hladkykh, Chris Poulin

DSG@softserveinc.com

v.02b-1

Contents

1	Introduction	2
2	Theory	2
2.1	Reinforcement Learning	2
2.2	Neural Networks	3
2.2.1	Architecture	3
2.2.2	Training	4
2.3	Deep Learning	4
2.3.1	Convolutional Neural Networks	5
3	Implementation	6
3.1	Process description	6
3.1.1	Explanation	6
3.1.2	Prediction	6
3.1.3	Algorithm	7
3.2	Wrappers	8
3.2.1	List of the supported wrappers	8
3.3	TNNF Framework	8
3.4	GPU Computing	9
3.5	MDP Environment	9
3.6	Preprocessing	9
3.7	Neural Network	9
3.7.1	Network topology	9
3.7.2	Network model	11
3.7.3	Training	11
3.8	Further improvements	12
4	Post-study section	13
4.1	Post-study analysis	13
4.1.1	Input Data	13
4.1.2	Analysis stages	13
4.1.3	Score assessment	14
4.1.4	Score differences assessment	14
4.1.5	Player efficiency assessment	16
5	Getting Started with DRL	18
5.1	Requirements	19
5.2	Virtual Machine	20
5.3	Installation	20
5.3.1	Needed packages	20
5.3.2	How to install	22
5.3.3	CPU mode	24
5.4	Launching and using	25
5.4.1	First launch	25
5.4.2	Project structure	26
5.4.3	Using existing models	27
5.4.4	Manual play mode	30
6	Experiments	30
7	Annex A1	34

1 Introduction

Although games are usually undertaken for enjoyment, they are also used as an educational tool or as a tool for simulating complex life situations for training in critical thinking and decision making. Now there are a great variety of games ranging from competitive sports, to board and video games, or even mind training focused games. Each game has its individual type, form, and structure.

While for much of their history, other humans were the logical opponent for games. But in early 1950s, when the Era of the Personal Computer came, this restriction was removed. The computer has become number-one opponent. Since then there was an enormous amount of competitions which tried to estimate a human brain's abilities in playing games in comparison to the state of the art computer algorithms.^{1 2} Although there was a great progress in developing such 'opponent agents' during the last years, the main goal of our project lies in the artificial intelligence of the 'human player simulating' game playing agent.

The DeepMind³ team had recently proposed to use Convolution Neural Networks to extract features from raw data (without any hand written features). This new 'unsupervised learning' approach provides a huge flexibility and algorithm abilities increase in order to work on different data and solve different problems, with noteworthy performance.

With credit to the prior projects by Deepmind⁴, our Data Science Group has also decided to develop a solution which would be able to achieve at or better performance in playing video games with computer than an expert human player. Our main goal is to completely generalize this idea and make it able to be applied to any video game. This particular use case could become a breakthrough in the study of Artificial Intelligence.

2 Theory

At this time, we are not about to introduce any new mathematical ideas or concepts in our project. Rather we want to put together currently existing state-of-the-art approaches and techniques to achieve a higher goal by combining them. In this chapter we will briefly explain the main underlying concepts of our solution such as Reinforcement Learning and Neural Networks, to let you understand 'how it actually works'. [Note: There will be a great deal of math and other process descriptions, so feel free to skip this section and jump directly to Implementation (or even Getting Started with DRL) if you're already familiar with these concepts.]

2.1 Reinforcement Learning

Reinforcement Learning (RL)⁵ is the set of algorithms (or approaches) inspired by behaviorist psychology. The idea behind it is that a computer algorithm (hereafter, agent) interacts with some unknown or semi-known environment (may be stochastic) using predefined and limited set of actions in a way that maximize some cumulative reward. The environment is typically formulated in terms of Markov Decision Process (MDP)⁶. It defines an environment ε as a two finite sets:

- $S = \{s_1, \dots, s_n\}$ – set of environment internal states;
- $A = \{a_1, \dots, a_k\}$ – set of environment legal actions;

¹[http://en.wikipedia.org/wiki/Watson_\(computer\)#Jeopardy.21](http://en.wikipedia.org/wiki/Watson_(computer)#Jeopardy.21)

²http://en.wikipedia.org/wiki/Deep_Blue_versus_Garry_Kasparov

³<http://deepmind.com/>

⁴<http://www.cs.toronto.edu/~vmnih/docs/dqn.pdf>

⁵http://en.wikipedia.org/wiki/Reinforcement_learning

⁶http://en.wikipedia.org/wiki/Markov_decision_process

- $\delta : S \times A \rightarrow S$ – transition function (usually unknown);

An agent interacts with environment ε in discrete time steps. On each step t in a state s_t it chooses an action $a_t \in A$ which is subsequently sent as a command to the environment. The environment changes agent's state from s_t to s_{t+1} . As a result of transition $[s_t, a_t, s_{t+1}]$ (means: from state s_t to state s_{t+1} by taking action a_t) agent immediately receives a reward r_{t+1} . The reward given by environment often is discounted by some so-called discount-factor $0 \leq \gamma \leq 1$. Thus we can define an expected discounted return after game termination by:

$$R = \sum_{t=0}^N \gamma^t r_{t+1} \quad (1)$$

here, N is the game termination time.

The agent makes its decisions on action to be taken in order to maximize an expected reward by following some action-value function which calculates quality of a state-action pair with respect to the past experience:

$$Q : S \times A \rightarrow R \quad (2)$$

Based on Bellman equation of optimality, we can define the optimal action-value function $Q^*(s, a)$ as the maximum expected reward for the state-action pair (s, a) achievable by taking action a in state s and thereafter following an optimal strategy by maximizing the expected value of $r + \gamma Q^*(s', a')$:

$$Q^*(s, a) = \mathbb{E}_{s' \sim \varepsilon} [r + \gamma \max_{a'} Q^*(s', a') | s, a] \quad (3)$$

The problem of optimal strategy determination, based on formulated optimization criterion, is referred to the dynamic programming problems. There are two general approaches:

- Iteration by strategy
- Iteration by values

Both approaches are computer-intensive, especially when power of possible states and action sets is significant enough. It explains inappropriateness of using traditional algorithms, when Markov Decision Process is reduced to hidden Markov model. In this case we can apply some approximate approaches, and one of the most successful solutions is applying Q-learning approximation by using apparatus of artificial neural networks.

The agent follows a greedy strategy on action selection $a = \max_a Q(s, a)$. However it can randomize its selection to explore states that have not been visited yet. Thus we define the ϵ -greedy strategy that follows the greedy strategy with probability $1 - \epsilon$ and selects a random action with probability ϵ .

2.2 Neural Networks

Artificial Neural Networks (ANNs) are mathematical models inspired by biological neural networks and are limited imitation of human's central nervous system (in particular the brain). In Machine Learning⁷ they are used as an approach to approximate functions $f : X \rightarrow Y$ (usually non-linear) by "learning" it from large amount of input data.

2.2.1 Architecture

There is a variety of approaches to build Neural Network topology, but in our project we consider only feedforward approach⁸. It contains following basic elements:

- Nodes (known as neurons) connected together by directed links to form a network;
- Connections between neurons associated with weights;
- Layers composed of neurons;

⁷http://en.wikipedia.org/wiki/Machine_learning

⁸http://en.wikipedia.org/wiki/Feedforward_neural_network

A connection from neuron i to neuron j propagates the activation a from i to j . Each connection also has a numeric weight $w_{i,j}$ associated with it, which determines the strength of the link. Each neuron j , in order to calculate its activation value a_j , first computes a weighted sum of all of its inputs, and then applies an activation function g to it:

$$a_j = g\left(\sum_{i=0}^n w_{i,j}a_i\right) \quad (4)$$

The activation function could be:

- sigmoid function ⁹.
- linear function ¹⁰.
- other advanced techniques (such as SoftMax, MaxOut etc)

Neural Networks usually are arranged in multiple layers of nodes in a directed graph, with each layer fully connected to the next one. The first one usually is called as Input layer and the final layer the Output layer. Intermediate layers of nodes between the Input and Output layers called the Hidden layer. The activation values on input layer represent the data input and the output layer activations represent the result value. The hidden layers are responsible for intermediate calculations that fit to the underlying data structure. The number of neurons in each layer depends on the amount of input data and the type (or complexity) of problem being solved.

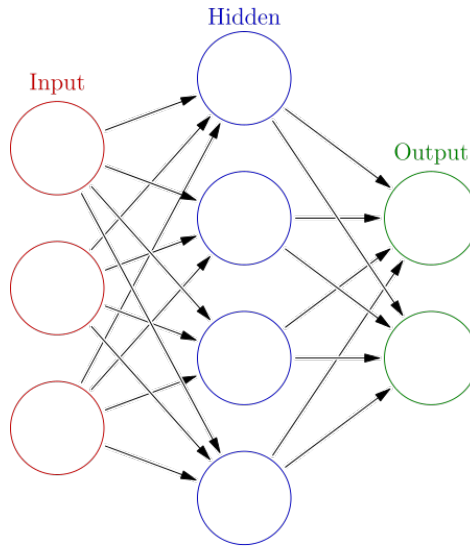


Figure 1: Neural network

2.2.2 Training

In Neural Networks learning occurs in changing connection weights by minimizing the Loss function. This is an example of supervised learning, and is carried out through backpropagation technique ¹¹ in conjunction with an optimization method such as gradient descent.

2.3 Deep Learning

Deep learning is an area of machine learning which is composed of a set of algorithms and techniques that attempt to define the underlying dependencies in a data and to model its high-level abstractions. The main goal of this approach is to avoid manual description of a data structure (like hand-written features) by automatic learning it from the data. Its name refers to the fact that typically any Neural Network with two or more hidden layers is called deep neural network. Deep Learning techniques usually

⁹http://en.wikipedia.org/wiki/Sigmoid_function

¹⁰http://en.wikipedia.org/wiki/Linear_function

¹¹<http://en.wikipedia.org/wiki/Backpropagation>

are applied to the tasks on a data which is not easy to represent as a valuable set of features because of its complexity (such as image or audio signal). Thus its main field of application lies in computer vision and speech recognition.

2.3.1 Convolutional Neural Networks

Convolutional neural network is a deep learning technique which is now widely used in image recognition. Architecturally, it is a type of feed-forward artificial neural network where the individual neurons are tiled in such a way that they respond to overlapping regions in the visual field.¹²

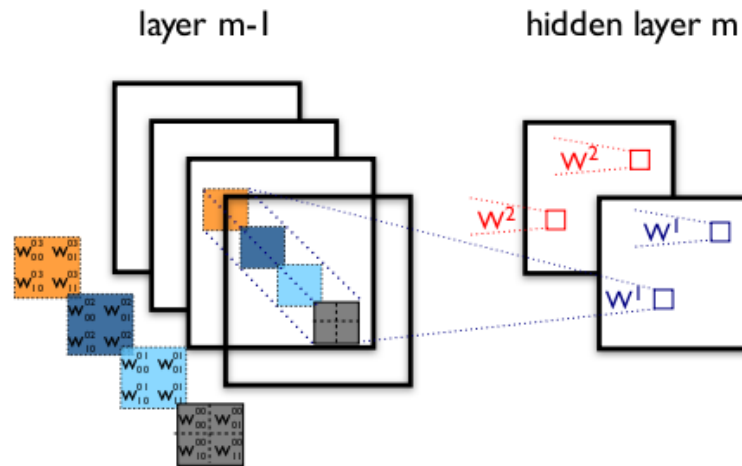


Figure 2: Convolutional neural network

Convolutional neural networks consists of a combination of multiple hidden layers of small cells of neurons which look specifically at small frame of the input image. The results of these cells are then tiled so that they overlap to obtain a better representation of the original image; this is repeated for every such layer. Because of this, they are able to tolerate translation of the input image.¹³ These cells act as local filters over the input space and are well-suited to exploit the strong spatially local correlation present in natural images.

Convolutional neural networks also implement a few important concepts (such as Shared Weights¹⁴ and MaxPooling¹⁵) that bring a big advantage in a field of Deep Learning and Computer Vision.

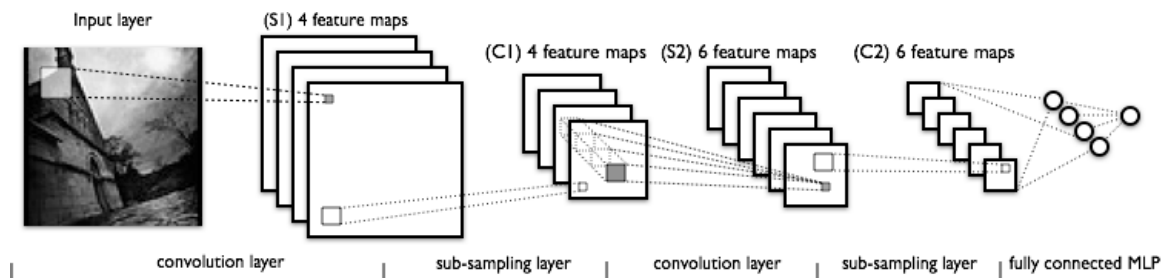


Figure 3: Full-stack model example

¹²http://en.wikipedia.org/wiki/Convolutional_neural_network#cite_note-deeplearning-1

¹³<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.125.3812>

¹⁴<http://deeplearning.net/tutorial/lenet.html#shared-weights>

¹⁵<http://deeplearning.net/tutorial/lenet.html#maxpooling>

3 Implementation

3.1 Process description

3.1.1 Explanation

The Diagram below (Figure 4) illustrates the flow of ‘train and predict’ functionality of the neural network. The blocks, marked by blue color and numbering are suggestions for future improvement. Meanwhile, there are also explanatory notes with benefits of each, with links for corresponding papers in chapter 3.8.

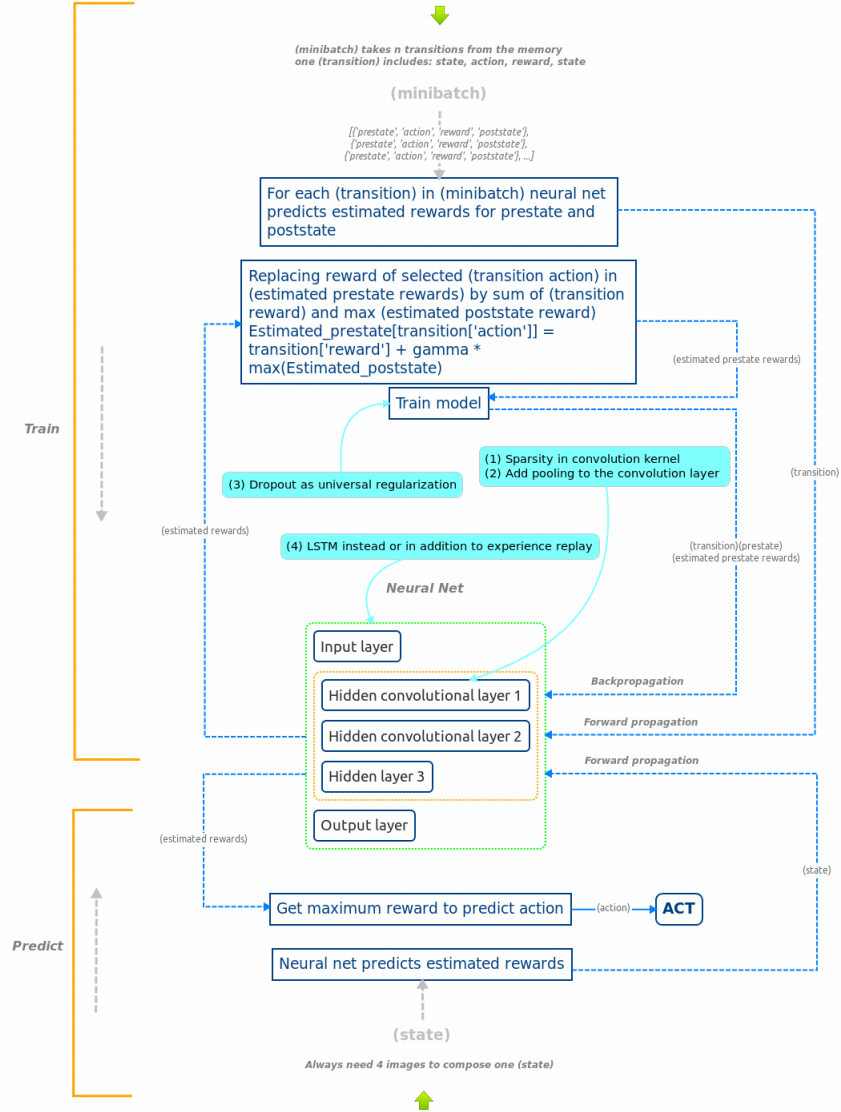


Figure 4: Learning diagram

3.1.2 Prediction

See Figure 4 where the network is designed to predict future rewards for all possible actions. There is a number of neurons that correspond to number of possible action in a game on the output layer (For Breakout there are four possible actions: no action, up (it's odd but possible), right, left).

- To make prediction we forward propagate the net with data describing current state at input and get activations of output layer.
- Then we use a greedy policy: define the neuron with maximum activation and perform corresponding actions.

For example, if the third neuron on output layer has a maximum activation we perform the third action. In diagram (Figure 4) this part of functionality starts from bottom to top. First block after bottom - forward propagation. And that one above this – illustrates a greedy policy algorithm.

3.1.3 Algorithm

Following algorithm was chosen for implementation in code (Figure 5).

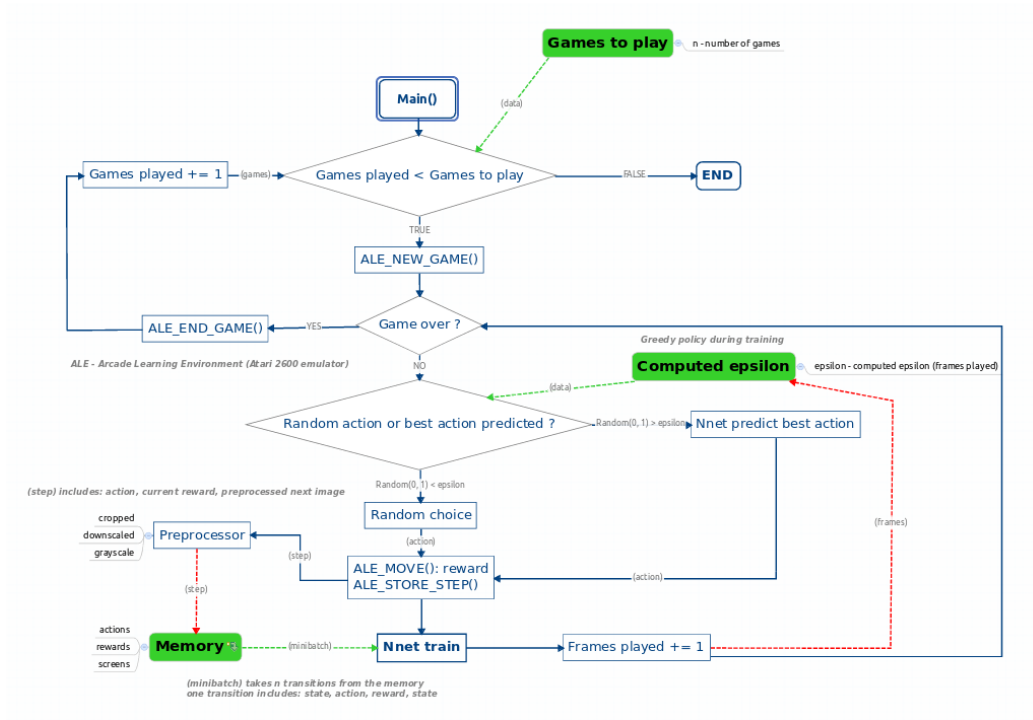


Figure 5: Implementation algorithm

Implementation of the memory model is relatively simple (Figure 6).

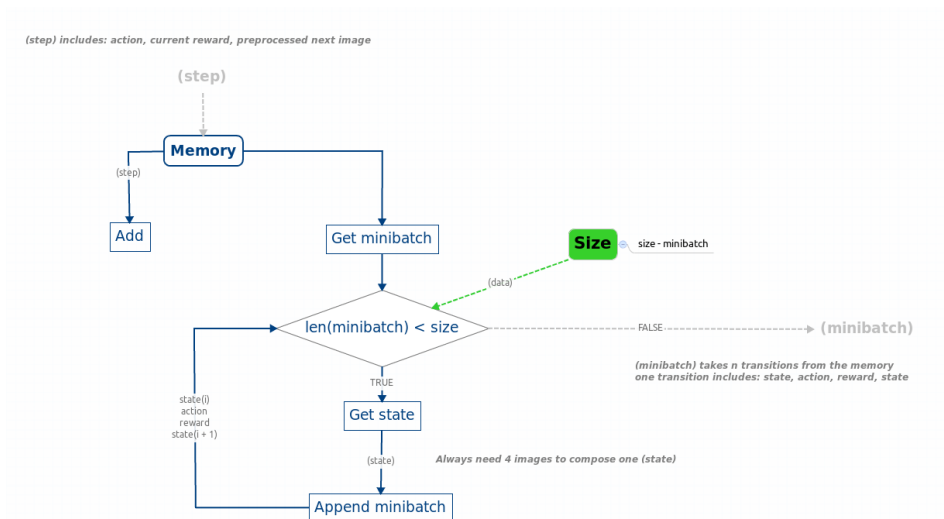


Figure 6: Memory model

3.2 Wrappers

Our library (DRL) provides a set of wrappers, each of these connects to appropriate game environment (Atari, Nintendo and others) using available software interfaces. The wrapper interacts with environment by sending it a commands to perform an actions or to obtain the necessary information about environment internal state (such as game reward, termination state etc). It also receives a video output from environment and processes it sequentially as an image stream. The wrapper is the only game-specific and non-generalized component in our solution. General diagram of DSG technology stack is provided on Figure 7.

3.2.1 List of the supported wrappers

Name	Type	Interface	Python class	Environments	Description
atari	Emulator	FIFO (text file)	ale.ALE	<ul style="list-style-type: none">• Breakout• Space Invaders	The Arcade Learning Environment (ALE) is a simple object-oriented framework that allows to develop AI agents for Atari 2600 games ¹⁶ .

3.3 TNNF Framework

Our Deep Learning implementation is based on TNNF framework which is also developed by our group. Please refer to the project documentation page ¹⁷ to get further information.

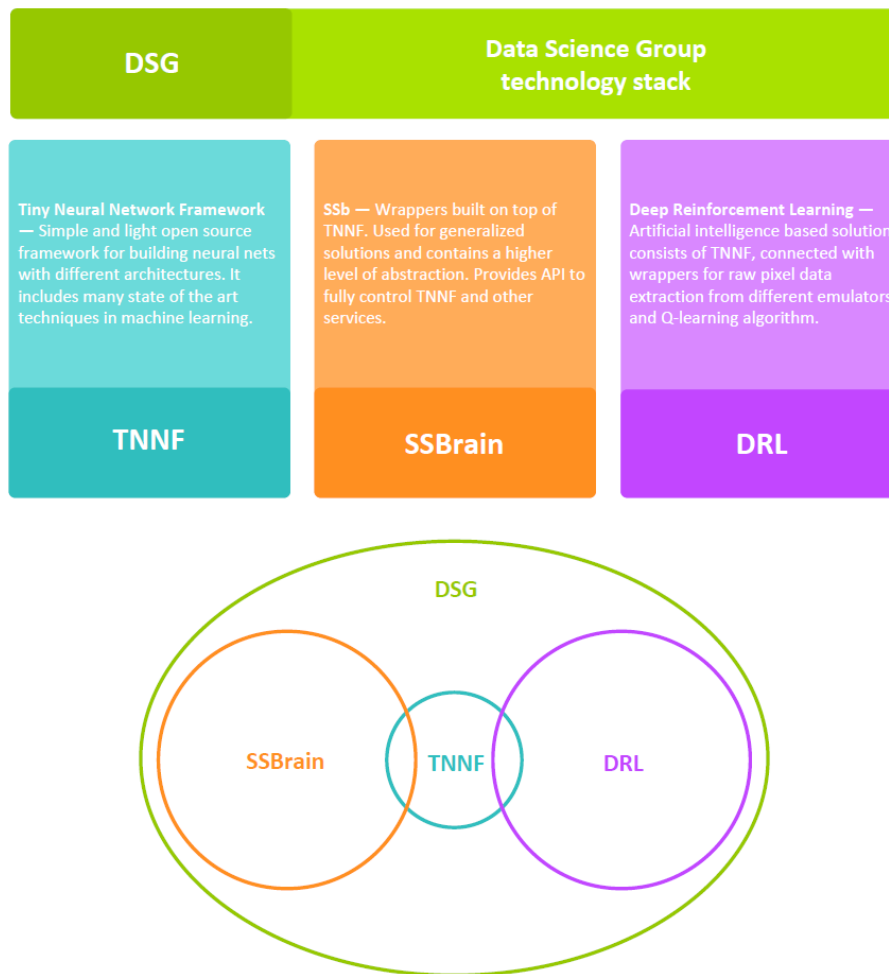


Figure 7: DSG technology stack

¹⁶<http://www.arcadelearningenvironment.org/>

¹⁷<http://tnnf.readthedocs.org/en/latest>

3.4 GPU Computing

Great progress achieved in last 10 years in the area of general-purpose computing on GPUs^{18 19} also brought a great benefit to the field of Artificial Intelligence (and to the Deep Learning in particular²⁰) since its incredible ability to massively parallel computation provides an exponential increase in performance (up to 50 times faster) on linear algebra operations that are so resource intensive with traditional CPU-based approach, especially on a large datasets.

There are two main frameworks that implements this computational model:

- NVIDIA CUDA (proprietary)²¹
- OpenCL (open source)²²

The neural network framework (TNNF) we use in this project, is built on the top of the Theano library²³ which in turn uses CUDA to perform GPU-based computations. Please refer to the project documentation²⁴ for more information on how to install and configure CUDA-backend in case if your GPU device is supported by vendor²⁵.

Our solution uses all advantages of this approach, however it still could be used in a traditional way of CPU computing (see Configuration)

3.5 MDP Environment

In our project, we consider the following MDP environment ε :

- Raw screen images from the game video output represent $s \in S$;
- Available game actions obviously represent set $A = \{a_1, \dots, a_k\}$;
- Game reward obtained through API represents r_t ;

3.6 Preprocessing

We perform a basic preprocessing step on each of raw 210×160 RGB images from input video stream. It consists:

- cropping a square 160×160 region of original image;
- downscaling to 84×84 image;
- converting RGB to gray-scale or B/W representation;

3.7 Neural Network

We use a Neural Network in conjunction with advanced Deep Learning techniques as a action-value function approximator. We refer to it as a DRL-network with weights θ which estimates $Q(s, a; \theta)$.

3.7.1 Network topology

DRL-network architecture consists of the following layers:

¹⁸http://en.wikipedia.org/wiki/General-purpose_computing_on_graphics_processing_units

¹⁹<http://www.nvidia.com/object/what-is-gpu-computing.html>

²⁰<http://devblogs.nvidia.com/parallelforall/accelerate-machine-learning-cudnn-deep-neural-network-library/>

²¹<http://en.wikipedia.org/wiki/CUDA>

²²<http://en.wikipedia.org/wiki/OpenCL>

²³<http://deeplearning.net/software/theano/>

²⁴http://deeplearning.net/software/theano/tutorial/using_gpu.html

²⁵http://en.wikipedia.org/wiki/CUDA#Supported_GPUs

1. Input layer (L1)

- Layer type: Convolutional
- Activation function: Leaky Rectified Linear Unit ²⁶
- Number of inputs: 28224
- Number of outputs: 6400
- Convolutional filter shape: $16 \times 4 \times 8 \times 8$

2. Hidden layer (L2)

- Layer type: Convolutional
- Activation function: Leaky Rectified Linear Unit
- Number of inputs: 6400
- Number of outputs: 2592
- Convolutional filter shape: $32 \times 16 \times 4 \times 4$

3. Hidden layer (L3)

- Layer type: Fully connected
- Activation function: Rectified Linear Unit ²⁷
- Number of inputs: 2592
- Number of outputs: 256

4. Output layer (L4)

- Layer type: Fully connected
- Activation function: Linear function
- Number of inputs: 256
- Number of outputs: *depends on a number of actions in particular game environment*

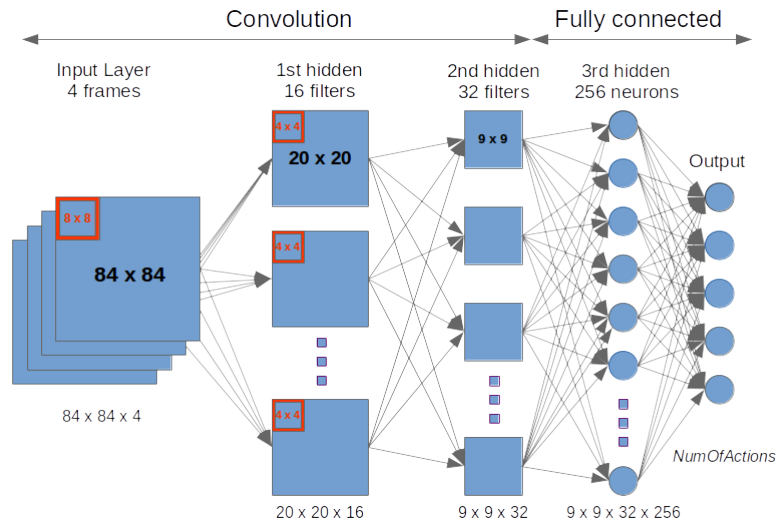


Figure 8: DRL Neural Network

²⁶[http://en.wikipedia.org/wiki/Rectifier_\(neural_networks\)#Leaky_ReLUs](http://en.wikipedia.org/wiki/Rectifier_(neural_networks)#Leaky_ReLUs)

²⁷[http://en.wikipedia.org/wiki/Rectifier_\(neural_networks\)](http://en.wikipedia.org/wiki/Rectifier_(neural_networks))

3.7.2 Network model

The following functions are used:

$$C'_p(i, j) = \ln(1 + \exp(\sum_{z=-3}^4 \sum_{v=-3}^4 (h'_{p_{zv}} \times S_u(2(i-3) + 1 + z, 2(j-3) + 1 + z)))) \quad (5)$$

$$C''_t(i, j) = \ln(1 + \max(\sum_{z=-1}^2 \sum_{v=-1}^2 (h''_{t_{zv}} \times C'_p(2(i-1) + 1 + z, 2(j-1) + 1 + z)), 0)) \quad (6)$$

$$F_k = \ln(1 + \max(\sum_{t=1}^{32} \sum_{i=1}^9 \sum_{j=1}^9 (w_{tijk}^C \times C''_t(i, j)), 0)) \quad (7)$$

$$Q_r = \sum_{k=1}^{256} (w_{kr}^F \times F_k) \quad (8)$$

where

$$t = (p-1) \times 2 + q_2; \quad p = (u-1) \times 4 + q_1; \quad q_1 = \overline{1, 4}; \quad q_2 = \overline{1, 2}; \quad u = \overline{1, 4} \quad (9)$$

S — input layer images

C' — first convolutional layer images

C'' — second convolutional layer images

F — hidden fully connected layer

Q — output layer

3.7.3 Training

A DRL-network is trained by minimizing a loss functions $L_i(\theta_i)$ on each iteration i , using RMSProp algorithm with mini-batches of size 128 samples:

$$L_i(\theta_i) = (y_i - Q(s, a, \theta_i))^2 \quad (10)$$

Here, $y = \mathbb{E}_{s' \sim \varepsilon}[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$ is the target return for iteration i .

Differentiating the loss function with respect to the weights we get the following gradient:

$$\nabla_{\theta_i} L_i(\theta_i) = \left[r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) - Q(s, a, \theta_i) \right] \nabla_{\theta_i} Q(s, a, \theta_i) \quad (11)$$

In terms of MDP, on each time-step t in states s_t, s'_{t+1} we divide environment timeline (Figure 9) into following parts:

- Pre-state - set of the last 4 image before state s_t ;
- Post-state - set of the last 4 image before state s'_{t+1} , when an action a_t has already been applied;

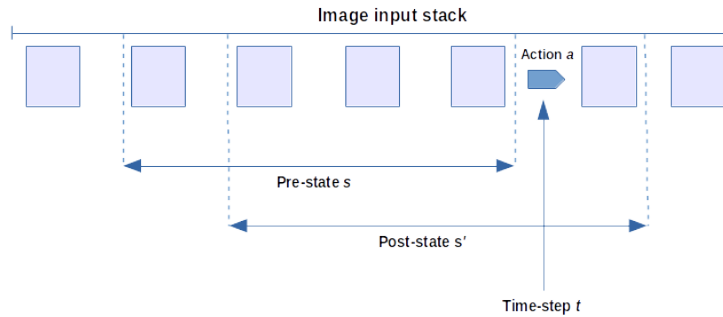


Figure 9: Timeline illustration

On each game iteration i , agent receives from environment wrapper a pre-state sequence s . Then it selects an action to be performed by following an ϵ -greedy strategy and using DRL-network to estimate $\max_a Q(s, a; \theta_i)$. Subsequently agent moves into post-state s' and gets a reward r . Reward takes value from set $\{-1, 0, 1\}$, where

- 1 - success (ball broke the brick)
- 0 - no action is required
- -1 - failure (ball was lost)

The sequence (s, a, r, s') forms a transition unit and is used to train DRL-network by estimating $y = \mathbb{E}_{s' \sim \varepsilon} [r + \gamma \max_{a'} Q(s', a'; \theta_{i-1}) | s, a]$. Complete training process workflow is reflected on Figure 10.

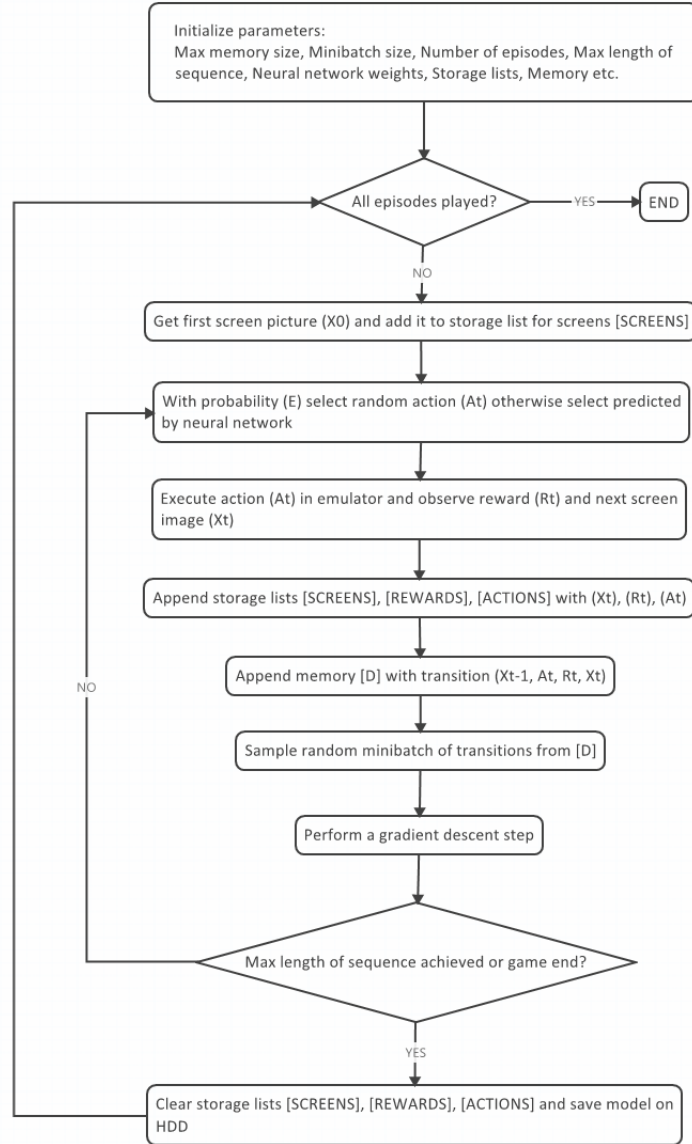


Figure 10: Complete workflow

3.8 Further improvements

There still exist, many possible options for our algorithm's improvement:

1. Initialization of connection weights to C-layers, based on preliminary analysis of the controlled object. We can assume that the convolution kernels, which were got after neural network training finishing, reduce to the set of "lines" and "spots" filters, so if the initialization was performed based on pre-set filters, it could significantly decrease required training time.

2. Network architecture can be extended by recurrent layer, for providing the possibility to get clustering of video series with unlimited number of clusters. This network can be trained to recognize the new information without retraining. This approach allows to generalize input information and use the optimal solution in the case, when we get new image (video series), but which is referred to the known class.
3. Sparsity helps extract better features from raw data, more complex and generalized. ²⁸
4. After obtaining features using convolution, we would next like to use them for classification. In theory, one could use all the extracted features with a classifier such as a softmax classifier, but this can be computationally challenging. To address this we could use aggregation operation called pooling. ²⁹
5. Dropout as universal regularization works better than weight decay and fix the overfitting problem. ³⁰
6. Long short term memory (LSTM), which is a recurrent neural network (RNN) architecture, can be used instead of experience replay technique. ³¹

4 Post-study section

It's a well-known fact that if we describe the target of research in view of Markov Decision Process, we can face the problem that the number of discrete states is huge, e.g. tens of thousands. Even though Q-learning algorithm allows to converge to the optimal value function, it might require a lot of time to find the real optimal strategy. The situation is further complicated by the necessity to determine the stop criterion and boundedness of learning strategy, so we can't guaranty that human couldn't find a more optimal solution. Especially, if we take into account the fact that once trained, network will not be adapted to the new conditions (reactions of controlled system). Human always adapts to the new conditions and learns during all his life.

So there is a hypothesis that the optimal results of real gamer (gamers) after statistically significant number of rounds will be higher than the trained neural networks. For the hypothesis checking we can gather both trained neural network and human game results and estimate the statistical significance of divergence between them (by using U-test for independent samplings).

Because this analysis requires significant statistical sampling, we, first of all, have to provide data gathering. For this purpose an additional module, which is made as "Human play manually" interface was developed.

4.1 Post-study Analysis

4.1.1 Input Data

Results of 30000 game iterations, represented in view of average core per one game, which were played during 1000 iteration.

4.1.2 Analysis stages

1. Score assessment
2. Score differences assessment
3. Player efficiency assessment

²⁸http://ufldl.stanford.edu/wiki/index.php/Autoencoders_and_Sparsity

²⁹<http://ufldl.stanford.edu/wiki/index.php/Pooling>

³⁰<http://jmlr.org/papers/volume15/srivastava14a/srivastava14a.pdf>

³¹http://en.wikipedia.org/wiki/Long_short_term_memory

4.1.3 Score assessment

For assessment of game score confidence level of the result we can estimate the result variance based on 95% confidence interval of fixed scores:

$$Score_{ConfL} = (1 - \frac{|ConfInt_{upper} - ConfInt_{low}|}{|Score_{est}|}) \times 100\% \quad (12)$$

$$Score_{est} = \frac{1}{n} \sum_{k=1}^n score_i \quad (13)$$

$$ConfInt_{upper} = Score_{est} + t_{0.975, n-1} (\frac{DScore_{est}}{\sqrt{n}}) \quad (14)$$

$$ConfInt_{low} = Score_{est} - t_{0.975, n-1} (\frac{DScore_{est}}{\sqrt{n}}) \quad (15)$$

$$DScore_{est} = \sqrt{\frac{1}{n-1} \sum_{k=1}^n (score_i - Score_{est})^2} \quad (16)$$

There are three diapasons of confidence level value, which can be considered as a satisfactory result:

- more than 80%, high statistical significance of estimated score
- from 60% to 80%, good statistical significance of estimated score
- from 50% to 60%, satisfactory statistical significance of estimated score

According to provided data, the following results can be observed:

Game	Player	$Score_{est}$	$DScore_{est}$	$ConfInt_{low}$	$ConfInt_{upper}$	$Score_{ConfL}$
Breakout	DRL	18,9	7,6	16,37	21,56	86%
	Human	11,5	8,6	7,58	15,37	66%
	RND	0,15	0,44	0,01	0,23	0,05%
Space Invaders	DRL	338,03	114,12	298,21	377,85	86%
	Human	356,658	160,3838	279,355	433,960	78%
	RND	135,443	62,589	112,877	158,008	83%
Boxing	DRL	12,20	8,198	9,14	15,26	75%
	Human	-7,20	9,764	-11,77	-2,63	33%
	RND	-7,27	9,176	-11,34	-3,20	44%

The frequency histograms and box diagrams are provided in 7.

As we can see, for games "Breakout" and "Space Invaders" estimated scores can be considered as statistically significant (except RND result for Breakout, which is explained by proximity of RND score to zero). "Boxing" game requires higher number of cases for providing more robust estimations, but in general, received score assessments can be considered as statistically significant enough.

4.1.4 Score differences assessment

To estimate statistical significance of difference between game scores based on DRL, Human and RND (random movements) player results, we can perform the means of the samples comparison analysis. On the table below results of Kolmogorov-Smirnov test for normal distribution are displayed:

Breakout:

		DRL	Human	RND
N		36	21	62
Normal parameters ^{a,b}	Mean	18,96575	11,4757	0,1202
	Std. dev.	7,67011	8,56729	0,44611
	Absolute	0,171	0,239	0,133
Most extreme differences	Positive	0,171	0,239	0,133
	Negative	-0,112	-0,141	-0,066
Kolmogorov-Smirnov Z		1,025	1,095	1,044
Asymp. Sig. (2-tailed)		0,244	0,182	0,226

Space Invaders:

		DRL	Human	RND
N		34	19	32
Normal parameters ^{a,b}	Mean	338,029	356,658	135,442
	Std. dev.	114,117	160,3838	62,589
	Absolute	0,142	0,129	0,177
Most extreme differences	Positive	0,142	0,129	0,177
	Negative	-0,096	-0,128	-0,101
Kolmogorov-Smirnov Z		0,828	0,562	1,000
Asymp. Sig. (2-tailed)		0,499	0,910	0,270

Boxing:

		DRL	Human	RND
N		30	20	22
Normal parameters ^{a,b}	Mean	12,20	-7,20	-7,27
	Std. dev.	8,198	9,764	9,176
	Absolute	0,090	0,199	0,134
Most extreme differences	Positive	0,065	0,105	0,111
	Negative	-0,090	-0,199	-0,134
Kolmogorov-Smirnov Z		0,494	0,890	0,628
Asymp. Sig. (2-tailed)		0,967	0,407	0,826

As we can see, results of Kolmogorov-Smirnov test allow us to consider data distribution as similar to normal distribution, so for sample means comparison we can use Paired-Samples t-Test:

Breakout:

		Paired Differences					t	df	Sig. (2-tailed)
		Mean	Std. Dev.	Std. Error Mean	95% Confidence				
					Lower	Upper			
Pair 1	DRL-Human	7,41	12,73	2,78	1,61	13,20	2,67	20,0	0,015
Pair 2	DRL-RND	18,86	7,74	1,29	16,24	21,48	14,62	35,0	0,000
Pair 3	Human-RND	11,26	8,53	1,86	7,38	15,14	6,05	20,0	0,000

Space Invaders:

		Paired Differences					t	df	Sig. (2-tailed)
		Mean	Std. Dev.	Std. Error Mean	95% Confidence				
					Lower	Upper			
Pair 1	DRL-Human	-44,18	173,95	39,91	-128,02	39,67	-1,11	18,0	0,283
Pair 2	DRL-RND	202,14	144,24	25,50	150,13	254,14	7,93	31,0	0,000
Pair 3	Human-RND	210,71	169,94	38,99	128,80	292,62	5,40	18,0	0,000

Boxing:

		Paired Differences					t	df	Sig. (2-tailed)
		Mean	Std. Dev.	Std. Error Mean	95% Confidence				
					Lower	Upper			
Pair 1	DRL-Human	22,90	11,83	2,64	17,36	28,44	8,66	19,0	0,000
Pair 2	DRL-RND	22,14	7,83	1,67	18,66	25,61	13,25	21,0	0,000
Pair 3	Human-RND	-0,45	10,05	2,25	-5,15	4,25	-0,20	19,0	0,843

So, we get the following results:

1. **Breakout game** – samplings, which refer to DRL, Human and RND scores, can be considered as heterogeneous, so the difference between game results can be considered as statistically significant.
2. **Space Invaders game** – samplings, which refer to DRL and Human scores, can be acknowledged as homogeneous samplings, so we can conclude, that the efficiency of DRL is equal to real player.
3. **Boxing game** - Human play scores very close to RND ones, so results verification is required.

4.1.5 Player efficiency assessment

We can consider the relative divergence of player score results in comparison with average human productivity as player efficiency for each case:

$$Efficiency_i = \frac{score_i - Score_{human}}{Score_{human}} \times 100\% \quad (17)$$

To provide more robust estimations, centroid moving average filtration was applied to the efficiency assessment:

$$Efficiency'_i = \frac{1}{2r+1} \sum_{k=-r}^r Efficiency_{i+k} \quad (18)$$

The descriptive statistics for this metric for "Breakout" game are provided below:

		MA(eff_DRL,10,10)	MA(eff_Human, 5, 5)	MA(eff_RND, 10, 10)
N	Valid	26	17	52
	Missing	36	45	10
Mean		64,5114	0,0	-94,1751
Median		67,1413	15,3721	-94,1349
Std. Dev		15,24606	30,74585	1,09937
Skewness		-0,539	-0,045	-0,262
Std. Err. of Skewness		0,456	0,550	0,330
Kurtosis		-0,494	-1,489	-0,682
Std. Err. of Kurtosis		0,887	1,063	0,650
Percentiles	10	37,8664	-32,2607	-95,6520
	20	52,0241	-26,0830	-95,2525
	30	59,0643	-19,6768	-94,9630
	40	62,8599	-10,2331	-94,4359
	50	67,1413	15,3721	-94,1349
	60	72,5347	20,4467	-93,8331
	70	73,7206	33,1532	-93,4182
	80	76,0807	36,4689	-93,1582
	90	85,1601	46,5686	-92,6580

It is easy to see, that DRL shows the more effective results – 65% efficiency increase in comparison with average Human result. At the same time random action choice results in 95% efficiency decrease. Results of the other games are represented on the following table:

Game	DRL	RND
<i>Space Invaders</i>	-6,7%	-72%
<i>Boxing</i>	90%	4,2%

Final results can be represented in view of the diagram:

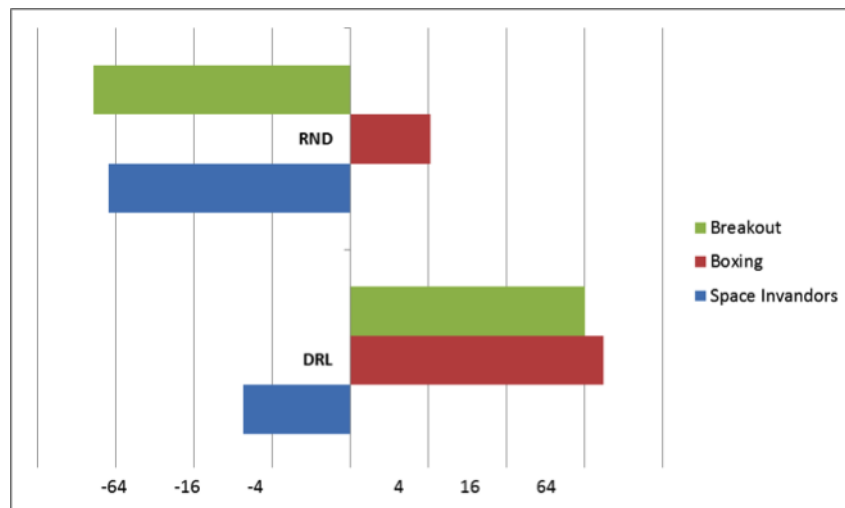


Figure 11: Final results

So, we can see that model allows to reach significant efficiency increase in comparison with random action selection, but the comparison of the results with human play requires the extension of the experiment and more detailed analysis.

5 Getting Started with DRL

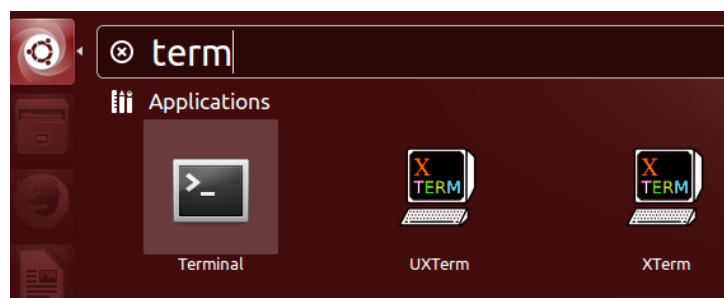
Questions, that are covered in this section:

- Installing all needed packages for project to run;
- Getting DRL source code and preparing for launching the training process;
- Structure of the project;
- Advance fitting.

All actions should be performed using the command line. Please, perform the next steps to launch terminal emulator in Ubuntu 14.04 or CentOS 7 (supported systems in this guide):

- Ubuntu 14.04

Press the **Alt-Ctrl-T** shortcut or use the graphic interface: press **Win**, then type “term” and choose your favorite:

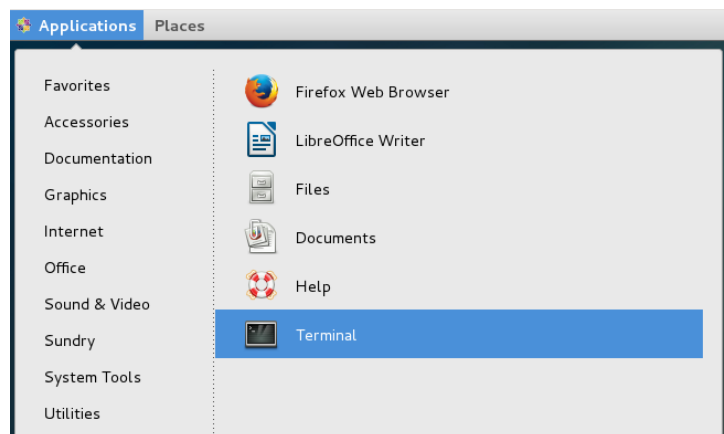


Starting terminal using GUI in Ubuntu

- CentOS 7

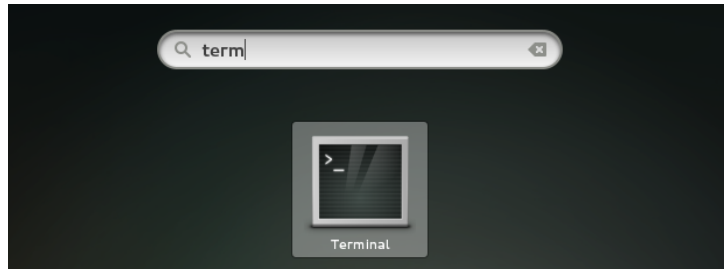
In CentOS there is no hotkey to start the terminal. But there is two ways to start it using GUI:

Using “Application” menu:



Starting terminal using GUI in CentOS with Application menu

Using **Win** button: press **Win**, then type “term” and choose your favourite



Starting terminal using GUI in CentOS with Win button

5.1 Requirements

Note

Please check the following requirements before you start an installation.

Supported operating systems (desktop versions):

- CentOS 7 ³²
- Ubuntu 14.04 ³³

Hardware requirements:

- **CPU:** 64-bit Intel or AMD multi-core processor (2+ Cores)
- **RAM:** 8GB of RAM (16GB recommended)
- **GPU:** CUDA supported video card ³⁴
- **Disk Space:** 8GB of free disk space

Make sure that your user is able to sudo. Run in terminal this command to do it:

```
sudo -l
```

The output should be like following:

Output

```
User <user> may run the following commands on <host>:
(ALL : ALL) ALL
```

```
drl@ubuntu: ~
drl@ubuntu:~$ sudo -l
[sudo] password for drl:
Matching Defaults entries for drl on ubuntu:
    env_reset, mail_badpass,
    secure_path=/usr/local/sbin\:/usr/local/bin\:/usr/sbin\:/usr/bin\:/sbin\:/bin
User drl may run the following commands on ubuntu:
    (ALL : ALL) ALL
drl@ubuntu:~$
```

Check sudo access

³²<http://www.centos.org/download>

³³<http://www.ubuntu.com/desktop>

³⁴<https://developer.nvidia.com/cuda-gpus>

If there is an error in the output, please ask your system administrator to provide you appropriate sudo rights.³⁵

5.2 Virtual Machine

For your convenience, we prepared a virtual machine image (in OVA format) containing fully operative demonstration setup of the DRL. With this image you would be able to test the pre-trained model, but without training, so you don't need CUDA supported video card and 4GB of RAM is enough.

Note

Demonstration build can be used only with a traditional CPU-based approach, since it is currently impossible to emulate a NVIDIA environment.

There are two distributions available in a Download section on a project web page:

- CentOS based
- Ubuntu Linux based

After extracting an image from archive, follow the instructions below:

1. Import an extracted VM image into your virtual environment (VMware Player, VirtualBox etc.);
2. Configure an instance according to the hardware requirements;
3. Start an instance;
4. Access to the virtual machine console and log into a desktop environment using following credentials:
 - Login: drl
 - Password: admin
5. Open the README file on the Desktop;
6. Follow the instructions in README.

Now you can test the pre-trained model. If you want to train a model by yourself, follow the instructions in the following sections.

5.3 Installation

Installation process was tested on supported systems (CentOS 7 and Ubuntu 14.04) installed on virtual machines from the official sites. Before installation there wasn't any additional software except for software, that was provided with this distributives.

In this section will be covered questions about what packages are needed for DRL to work properly and how to install them.

5.3.1 Needed packages

See "How to install" section (5.3.2) to perform installation of this libraries.

Names of packages are provided in format:

Conventions

<Ubuntu_package>/<CentOS_package> - names of packages in both supported systems delimited by slash (/)

OR

³⁵<http://www.pendrivelinux.com/how-to-add-a-user-to-the-sudoers-list/>

Conventions

<Name_of_package> (System) - if package is needed to be installed only

General purpose packages:

- *g++/gcc-c++* - the C++ front end for GCC compiler.
- *libopenblas-dev/blas-devel* - The BLAS (Basic Linear Algebra Subprograms) are routines that provide standard building blocks for performing basic vector and matrix operations. ³⁶
- *libyaml-dev/libyaml-dev* - YAML is a human friendly data serialization standard for all programming languages. Is needed for *pylearn2*. ³⁷
- SDL packages:

libsdl1.2-dev/SDL-devel.x86_64

libsdl-image1.2-dev/SDL_image-devel.x86_64

libsdl-gfx1.2-dev/SDL_gfx-devel.x86_64

SDL - Simple DirectMedia Layer is a cross-platform development library designed to provide low level access to audio, keyboard, mouse, joystick, and graphics hardware via OpenGL and Direct3D. It is used by video playback software, emulators, and popular games. ³⁸

- *gcc (CentOS)* - compiler system produced by the GNU Project supporting various programming languages. ³⁹
- *gcc-gfortran (CentOS)* - free Fortran 95/2003/2008 compiler for GCC. ⁴⁰
- *lapack-devel (CentOS)* - provides routines for solving systems of simultaneous linear equations, least-squares solutions of linear systems of equations, eigenvalue problems, and singular value problems. ⁴¹
- *atlas-devel (CentOS)* - provides C and Fortran77 interfaces to a portably efficient BLAS implementation. ⁴²
- *zlib-devel.x86_64 (CentOS)* - general purpose compression library. ⁴³

Python packages:

This section is about what packages should be installed before DRL installation. The common commands for installation of this packages you can find in the next section. This information is provided for you to be informed what libraries do DRL use. Some of this packages need to be installed using **pip**. To find out which tool to use - see the next section.

List of packages:

- *python-dev/python-devel* - contains the header files and libraries needed for python to be extended with dynamically loaded extensions. ⁴⁴
- *python-pip/pip* - recommended tool for installing Python packages. ⁴⁵

³⁶<http://www.netlib.org/blas/>

³⁷<http://www.yaml.org>

³⁸<https://www.libsdl.org/>

³⁹<https://gcc.gnu.org/>

⁴⁰<https://gcc.gnu.org/wiki/GFortran>

⁴¹<http://www.netlib.org/lapack/>

⁴²<http://math-atlas.sourceforge.net/>

⁴³<http://www.zlib.net/>

⁴⁴<https://www.python.org/>

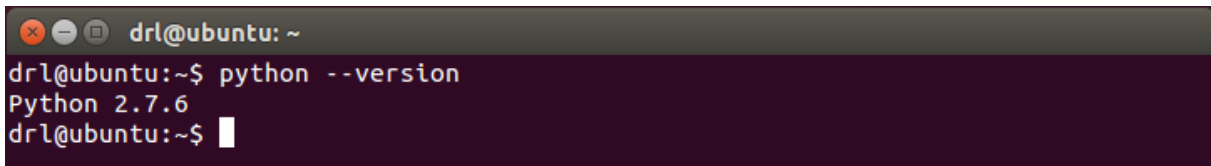
⁴⁵<https://pypi.python.org/pypi/pip/>

- *python-numpy/numpy==1.6.1* - the fundamental package for scientific computing with Python. ⁴⁶
- *python-scipy/scipy==0.10.1* - Python-based ecosystem of open-source software for mathematics, science, and engineering. ⁴⁷
- *python-matplotlib/python-matplotlib* - Python 2D plotting library. ⁴⁸
- *Theano/Theano* - a Python library that allows you to define, optimize, and evaluate mathematical expressions involving multi-dimensional arrays efficiently. ⁴⁹
- *pylearn2/pylearn2* - machine learning library. Needed in DRL for Convolution Networks. ⁵⁰
- *python-pil/python-pillow* - Python imaging library. ⁵¹
- *python-nose/python-nose* - testing library for Python. ⁵²

5.3.2 How to install

DRL is written in Python, so make sure that you have at Python 2.7 or later installed on your system. It may be easily checked from the command line:

```
python --version
```



```
drl@ubuntu: ~
drl@ubuntu:~$ python --version
Python 2.7.6
drl@ubuntu:~$
```

Check Python version

If the output looks like “*Python 2.7.x*” (Python 3.x.x is not ok) - you have the correct version of Python and you can skip Python installation section. In case of error or another version, please, install the correct version of Python on your system and, if needed, make a symlink “*python*” in one of the \$PATH directories to the interpreter of this version.

Python installation:

- Ubuntu 14.04:

```
sudo apt-get install build-essential
sudo apt-get install libreadline-gplv2-dev libncursesw5-dev libssl-dev \
libsqlite3-dev tk-dev libgdbm-dev libc6-dev libbz2-dev
cd ~/Downloads/
wget http://python.org/ftp/python/2.7.5/Python-2.7.5.tgz
tar -xvf Python-2.7.5.tgz
cd Python-2.7.5
./configure
make
sudo make install
```

⁴⁶<http://www.numpy.org/>

⁴⁷<http://www.scipy.org/>

⁴⁸<http://matplotlib.org/>

⁴⁹<http://deeplearning.net/software/theano/>

⁵⁰<https://github.com/lisa-lab/pylearn2>

⁵¹<http://www.pythonware.com/products/pil/>

⁵²<https://nose.readthedocs.org/en/latest/>

- CentOS 7:

```
sudo yum groupinstall "Development tools"
sudo yum install zlib-devel bzip2-devel openssl-devel ncurses-devel \
sqlite-devel readline-devel tk-devel gdbm-devel db4-devel libpcap-devel \
xz-devel
cd ~/Downloads/
wget --no-check-certificate \
https://www.python.org/ftp/python/2.7.5/Python-2.7.5.tar.xz
tar -xvf Python-2.7.5.tar.xz
cd Python-2.7.5
./configure --prefix=/usr/local
make
sudo make altinstall
```

Pre-installation:

If you are using CentOS, you should add *EPEL* and *RPMForge* repositories:

```
sudo yum install \
http://dl.fedoraproject.org/pub/epel/7/x86_64/e/epel-release-7-5.noarch.rpm
sudo yum install \
http://pkgs.repoforge.org/rpmforge-release/rpmforge-release-0.5.3-1.el7.rf.x86_64.rpm
```

Paste following lines line by line to your terminal to install all needed packages:

- Ubuntu 14.04:

```
sudo apt-get install python-pil python-numpy python-scipy python-dev \
python-pip python-nose g++ libopenblas-dev git libsdl1.2-dev \
libsdl-image1.2-dev libsdl-gfx1.2-dev python-matplotlib libyaml-dev
sudo pip install -U numpy
sudo pip install -U pillow==2.7.0
sudo pip install Theano
```

- CentOS 7:

```
sudo yum install python-devel python-nose python-pillow python-setuptools \
gcc gcc-gfortran gcc-c++ blas-devel lapack-devel atlas-devel \
SDL-devel.x86_64 SDL_image-devel.x86_64 SDL_gfx-devel.x86_64 \
zlib-devel.x86_64 git python-matplotlib libyaml-dev
sudo easy_install pip
sudo pip install -U numpy
sudo pip install -U pillow==2.7.0
sudo pip install scipy
sudo pip install Theano
```

Also, you should install *pylearn2* package. Steps are the same for the both systems:

```
cd ~/ # Or any other directory where you want the pylearn2 project to be stored
git clone https://github.com/lisa-lab/pylearn2.git
cd ./pylearn2
sudo python setup.py develop
```

DRL project is developed and tested with CUDA (parallel computing platform and programming model created by NVIDIA⁵³). CUDA is developed for fast GPU computing. CPU mode is also available,

⁵³<https://developer.nvidia.com/cuda-zone>

but deep reinforcement learning in this mode is too slow, so CPU mode is not recommended for learning. However, if you have a trained model, it is possible to test it on CPU without any speed issues.

Follow the instructions on NVIDIA page to install CUDA SDK⁵⁴.

Also, if you want to use CUDA not only in this terminal session, you have to change some global variables in your `~/.bashrc`:

```
cat >> ~/.bashrc <<EOF
export PATH=/usr/local/cuda-6.5/bin:$PATH
export LD_LIBRARY_PATH=/usr/local/cuda-6.5/lib64:$LD_LIBRARY_PATH
EOF
```

Great! Your system is ready for DRL to be installed and tested.

DRL installation:

To install DRL you need to:

1. Make sure that correct Python version is installed in your system and Python interpreter can be started with “python” command.
2. Make all the pre-installation steps. Make sure that everything is installed properly without errors.
3. Clone source code from GitHub⁵⁵:

```
git clone https://github.com/DSG-SoftServe/DRL
```

4. Navigate to the project directory

```
cd DRL
```

5. Run `install.sh` script.

```
bash ./install.sh
sudo chown -R $USER ~/.theano # to run the demo without sudo
```

If installation is successful, the environment is ready to use.

5.3.3 CPU mode

It is not recommended to run the demo in CPU mode. But, if you don't have NVIDIA GPU and want to check out if it's working, apply this fix:

```
cd src
drl_optimize=False # False for CPU, True for GPU
drl_regex="(\\s*optimized=)[a-zA-Z]+(\\))"
sed -ri "s:${drl_regex}\\1${drl_optimize}\\2:g" ai/neuralnet.py
```

WARNING

Use only "False" and "True" values

To launch the demo in CPU mode, use `cpuRun.sh` instead of `gpuRun.sh` in future paragraphs.

⁵⁴<http://docs.nvidia.com/cuda/cuda-getting-started-guide-for-linux/index.html#pre-installation-actions>

⁵⁵<https://github.com/DSG-SoftServe/DRL>

5.4 Launching and using

5.4.1 First launch

You have DRL installed and it's time for the first launch.

Navigate to the DRL project directory and perform following steps to launch the learning process:

```
cd ./src/  
bash gpuRun.sh main_exp.py
```

Congratulations, you have launched the deep reinforcement learning process.
Now you should see the window like this:



Figure 12: Breakout

Also, there should be an output in the command line like this:

Output

```
action RAND: 2  
S: 50 A: 49 R: 49 D: 171  
_____ Reward 0  
Step: 177  
0.04417
```

Let's find out what does this indicators means:

- *action RAND: 2* - there is 4 available actions ($[0..3]$). The action can be chosen from neural network (NN) experience with probability ϵ , where $0 < \epsilon < 1$. Than you'll see "NEURALNET" instead of "RAND". Otherwise, NN will perform a random action;
- *S: 50* - current state. States are changing with ascending order starting from 0 since last negative reward;
- *A: 49* - performed action. Logic of numbering is similar to state;
- *R: 49* - reward for action. Logic of numbering is similar to state;
- *D: 171* - memory;
- *Reward 0* - current value of reward. Could be -1, 0 or 1;

- *Step: 177* - number of steps performed from the learning start;
- *0.04417* - time, spent for the step (in minutes).

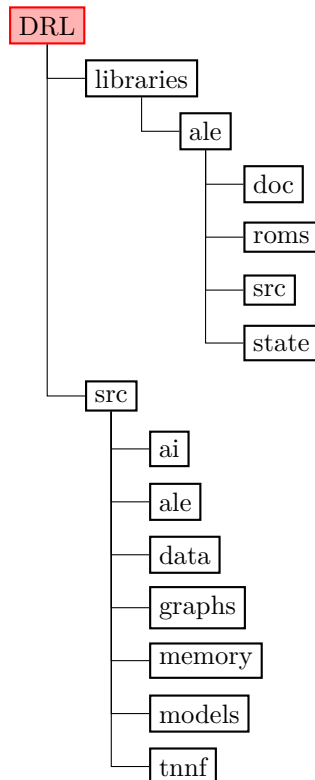
“S”, “A” and “R” are from the Markov Decision Process (MDP) model⁵⁶.

So, we can see such correlations between command line and window:

- when the “action” is changing, the actor (bottom block) changes it’s behavior;
- when the shell hits one of the upper blocks, the reward equals 1;
- when the actor loses the last life, the reward equals -1.

5.4.2 Project structure

Let’s look closer at some of the important parts of the project. Below, you can see the output of Unix “tree -d” command, except for unnecessary for understanding folders:



There is two folders in the root of the project structure that need to be mentioned:

- **libraries** - contains “ale” folder. ALE is a tool, which gave the ability to learn models on different Atari games. From the official site:
The Arcade Learning Environment (ALE) is a simple object-oriented framework that allows researchers and hobbyists to develop AI agents for Atari 2600 games. It is built on top of the Atari 2600 emulator Stella and separates the details of emulation from agent design.⁵⁷ The most interesting folder in libraries/ale directory is roms. Here you can find games, for which you can train NN. (see the additional fitting chapter)
- **src** - the main folder obviously. All source code of the DRL project is located in this folder. Five directories should be mentioned here:
 - ai* - source code of neural network, which is used within the project;
 - ale* - wrapper code for interaction with the Arcade Learning Environment;
 - data* - scripts to work with data;

⁵⁶http://en.wikipedia.org/wiki/Markov_decision_process

⁵⁷<http://www.arcadelearningenvironment.org/>

graphs - graphs with experiment results;
memory - implementation of neural network memory;
models - directory, where all trained models will be located;
tnnf - source code for Tiny Neural Net Framework. ⁵⁸ used to create and train the neural network, which is used in this project.

The **graphs** directory contains following graphs:

- *rewargs.png* - sum of rewards for multiple iterations (1000 by default), e.g. Left plots on Figures 15 and 16;
- *rgstat.png* - average rewards for 1 finished game from multiple iterations (1000 by default), e.g. Right plots on Figures 15 and 16;
- *wsum.png* - sum of the weights on both convolutional layers of neural network.

5.4.3 Using existing models

The result of the training part is a file, saved in *DRL/src/models/* directory. You can use this file to find out how good your NN was trained. You can use the *play.py* script in the *src* directory to play a game with your trained NN. Launching game with a trained NN is similar to training launching:

```
cd DRL/src
bash gpuRun.sh play.py
```

You should see the similar to the training output: window with the game and output in command line. But there shouldn't be any training time in the output.

In the first version of DRL there is no ways to configure which model to use, except for editing *play.py* script. By default, the pre-trained model “*~p-AS.model.weights.NN.NMem.906508*” is hardcoded.

You can change the script to use your own model in two ways: to edit the file or to run the *sed* tool.

Note

path to the model is relative to the *play.py* script

Edit the file using text editor

1. In command line navigate to the DRL project directory, then to *src* folder. Using your favorite text editor, open the *play.py* script;
2. Perform search to navigate to the line, which starts with “*self.nnet.loadModel*”;
3. Change the path in quotes to the path of your model;
4. Save the file.

```
19 # Loading model
20 self.nnet.loadModel('model.weights.NN.old.pillow')
21
```

String to change

⁵⁸<http://tnnf.readthedocs.org/en/latest/index.html>

Edit the file using sed

Also, you can use the sed command to change the file. Navigate to the DRL project directory, then to src folder and run this command to do it:

```
drl_model=./models/model.weights.NN.old.pillow # change this value
sed -ri "s:(\s*self\.nnet\.loadModel\(').*?(')):1$drl_model\2:" play.py
```

WARNING

Don't change the sed command except for the path to your model, unless you know what you are doing.

Additional fitting

In the first version of the DRL project you can't change training configuration easily using config files. But, if you really want to set up your own training parameters, you can do it by editing source files. Below you can see explanation of the most important training parameters and then, instruction how to change them.

Training parameters:

- *Memory size* - how many samples should we store in memory to take mini-batch from at one step;
- *Mini-batch size* - how many samples should we take from memory to perform the next learning step;
- *Decrease of epsilon step* - decreasing rate of epsilon which is responsible for decision to take a random step or a step, predicted by NN;
- *Number of allowed actions* - number of actions which NN can perform;
- *Skip frames* - there is too many frames which the emulator produces. Some of them should be skipped;
- *Maximum sequence length* - maximum number of steps that NN can perform for one game;
- *Maximum number of episodes* - maximum number of games during training;
- *Saving counter threshold* - when should the program save the game. Counts relatively to steps performed by NN;
- *Base model name* - the first part of the name of model which will be generated;

How to change parameters:

All this parameters are changing in file *DRL/src/main.py*, so navigate to the *src* directory:

```
cd DRL/src
```

In examples below will be used default values for you to be able to restore the initial state. Just change the value of the variable in the first command to change the certain training parameter.

WARNING

It is supposed that you are working with the unedited version of *"main.py"* of the first version of DRL. If you are not sure, please, don't use this commands. They may take no effect or even corrupt the code. Also, you may change this parameters in your favorite text editor, but it's not recommended.

Memory size:

```
drl_mem_size=100000 # change this value
drl_regex='(\s*self\.size\s?=\s?)[0-9]+(\s*#\s?(M|m)emory\s*size)'
sed -ri "s:$drl_regex:\1$drl_mem_size\2:" main_exp.py
```

Minibatch size:

```
drl_minibatch=128 # change this value
drl_regex='(\s*self\.minibatch_size\s?=\s?)[0-9]+(\s*#\s?(M|m)inibatch\s*size)'
sed -ri "s:$drl_regex:\1$drl_minibatch\2:" main_exp.py
```

Decrease of epsilon step:

```
drl_eps_step=700000 # change this value
drl_regex='(\s*self\.epsilon_step_dec\s*=\s*)([0-9]+)(\s*\.*)'
sed -ri "s|$drl_regex|\1$drl_eps_step\2|" main_exp.py
```

Number of allowed actions:

```
drl_actions=4 # change this value
drl_regex='(\s*self\.number_of_actions\s?=\s?)[0-9]+(\s*#\s?(L|l)eft,?\s*right.*)'
sed -ri "s:$drl_regex:\1$drl_actions\2:" main_exp.py
```

Skip frames:

```
drl_skipframes=6 # change this value
drl_regex='(\s*self\.ale\s*=\s*ALE\(\s*self\.memory,\s*skip_frames=\s?)[0-9]+(\s*\.*)'
sed -ri "s|$drl_regex|\1$drl_skipframes\2|" main_exp.py
```

Maximum sequence length:

```
drl_max_seqlength=50000 # change this value
drl_regex='(\s*while\sself\.T\s<\s)[0-9]+(\s*and\snot\sself\.ale\.game_over:)'
sed -ri "s|$drl_regex|\1$drl_max_seqlength\2|" main_exp.py
```

Maximum number of episodes:

```
drl_max_episodes=50000000 # change this value
drl_regex='(\s*while\sself\.M\s<\s)[0-9]+(:)'
sed -ri "s|$drl_regex|\1$drl_max_episodes\2|" main_exp.py
```

Saving counter threshold:

```
drl_count_threshold=20000 # change this value
drl_regex='(\s*if\sself\.count\s>=\s)[0-9]+(:)'
sed -ri "s|$drl_regex|\1$drl_count_threshold\2|" main_exp.py
```

Base model name:

```
drl_model_name=AS.model.weights.NN.NMem # change this value
drl_regex='(self\.nnet\.saveModel\(''\.\./models\/\).*('\s\+\sstr\(\s*self\.steps\)\))'
sed -ri "s|$drl_regex|\1$drl_model_name\2|" main_exp.py
```

To revert all changes:

```
drl_regex='(\s*self\.size\s?=\s?)[0-9]+(\s*\s?(M|m)emory\s*size)'
sed -ri "s:${drl_regex}:1100000\2:" main_exp.py
drl_regex='(\s*self\.minibatch_size\s?=\s?)[0-9]+(\s*\s?(M|m)inibatch\s*size)'
sed -ri "s:${drl_regex}:1128\2:" main_exp.py
drl_regex='(\s*self\.epsilon_step_dec\s*=\s*)[0-9]+(\s*\s?(L|l)eft,?\s*right.*)'
sed -ri "s:${drl_regex}:1700000\2|" main_exp.py
drl_regex='(\s*self\.number_of_actions\s?=\s?)[0-9]+(\s*\s?(L|l)eft,?\s*right.*)'
sed -ri "s:${drl_regex}:14\2:" main_exp.py
drl_regex='(\s*self\.ale\s*=\s*ALE\(\self\.memory,\sskip_frames=\)[0-9]+(\s*\s?(L|l)eft,?\s*right.*)'
sed -ri "s:${drl_regex}:16\2|" main_exp.py
drl_regex='(\s*while\sself\.M\s<\s)[0-9]+(:)'
sed -ri "s:${drl_regex}:15000000\2|" main_exp.py
drl_regex='(\s*while\sself\.T\s<\s)[0-9]+(\sand\snot\sself\.ale\.game_over:)'
sed -ri "s:${drl_regex}:150000\2|" main_exp.py
drl_regex='(\s*if\sself\.count\s>=\s)[0-9]+(:)'
sed -ri "s:${drl_regex}:120000\2|" main_exp.py
drl_regex="(self\.nnet\.saveModel\(\'\.\/models\/\).*\(\.\'\s+\sstr\(\self\.steps\)\)\)"
sed -ri "s:${drl_regex}:1AS.model.weights.NN.NMem\2|" main_exp.py
```

5.4.4 Manual play mode

In order to compare results of a DRL Agent and a human, we created an interface for a manual play in the current environment.

You have DRL Agent installed and trained so it's time for the manual play. Navigate to the DRL project directory and perform following steps to launch the playing process:

1. Go to the source directory

```
cd DRL/src
```

2. Run *human_play.py* script

```
bash cpuRun.sh human_play.sh
```

3. To control the game after launch select the terminal window and send commands from the keyboard

To figure out control buttons, open *human_play.py* script and navigate to the relevant section. You will find a list of actions corresponding to a keyboard buttons there (Figure 13). To change the game played:

1. Restart terminal
2. Edit *ale.py* script and change the games name as it is shown on Figure 14.

Note: restart the terminal after every game series played for further tasks performing.

6 Experiments

The model was trained with 900000 steps on NVIDIA GeForce GTX 880M. The training time was two days. This model was compared to random action model. The trained model, as expected, showed better results. The analysis consists of two comparisons: of total rewards for 1000 steps and of average games rewards for 1000 steps.

The experiment was evaluated on 100000 steps. For one hit, the agent gets reward equals 1. On graphs, one point - is result for 1000 steps. On the "Y" axes there is a reward and on the "X" axes there is a number on batch.

Below you can see raw results for the random action model (Figures 15):

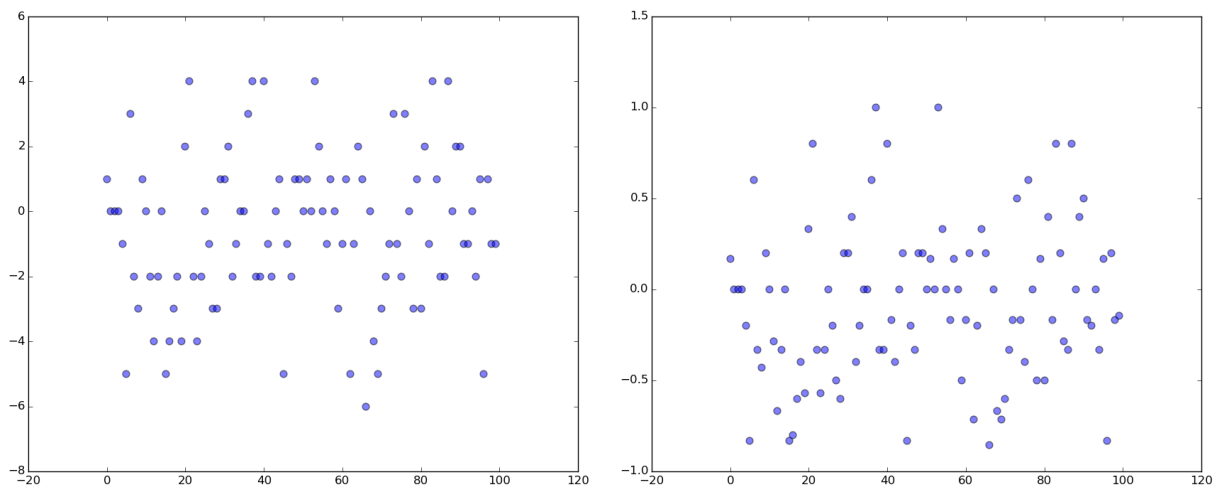
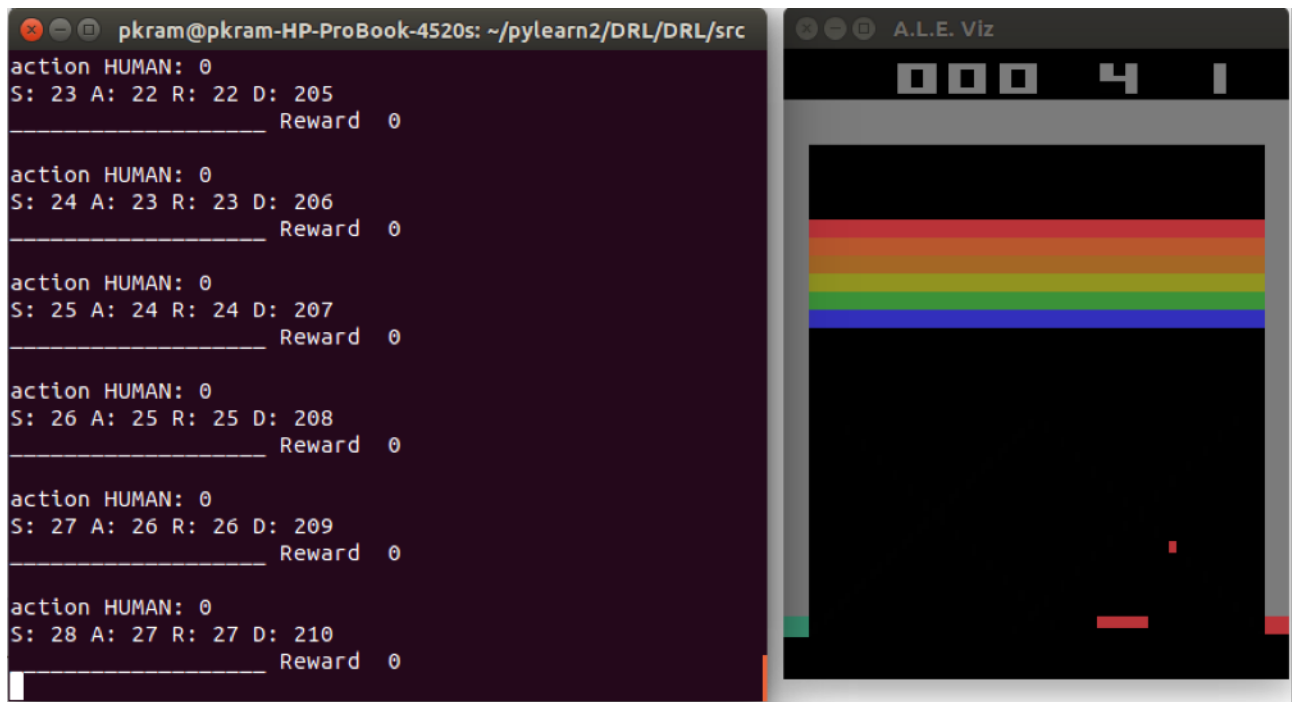


Figure 15: Random action model. Left - rewards for 1000 iterations; Right - average rewards for 1 finished game.

The graphs shows expected behavior. Mean rewards are around 0.
But for the trained model - rewards are higher (Figures 16):


```

human_play.py x
# Debug
#print "Key", char
if not(t.isAlive()):
    t = threading.Thread(target=keypress)
    t.daemon = True
    t.start()

#-----
if str(char) == "a":
    action = 2
elif str(char) == "w":
    action = 3
elif str(char) == "/":
    action = 1
elif str(char) == "d":
    action = 5
elif str(char) == "s":
    action = 4
else:
    action = 0
print '\naction HUMAN:', action

# We collect all moves, so using special methods
# we can save them to use for training process later:
# self.memory.mem_save(self.steps)

# Act and store action in memory

```

Figure 13: Game control bindings

```

#-----#
# Atari emulator wrapper.
#-----#
import os
import numpy as np
from time import sleep
from data.preprocessing import preprocessing
from TNNF import fTheanoNNclassCORE, fGraphBuilderCORE
#-----#

class ale:
    def __init__(self, memory, display_screen="true", frames_to_skip=4, ale_game_ROM='../emulators/ale_04/roms/breakout.bin'):
        # List of possible actions for agent
        self.actions_list = [np.uint8(0), np.uint8(1), np.uint8(2), np.uint8(3), np.uint8(4), np.uint8(5)]
        # Read and write commands and response from emulator
        self.f_in = ""

```

Figure 14: Game selection

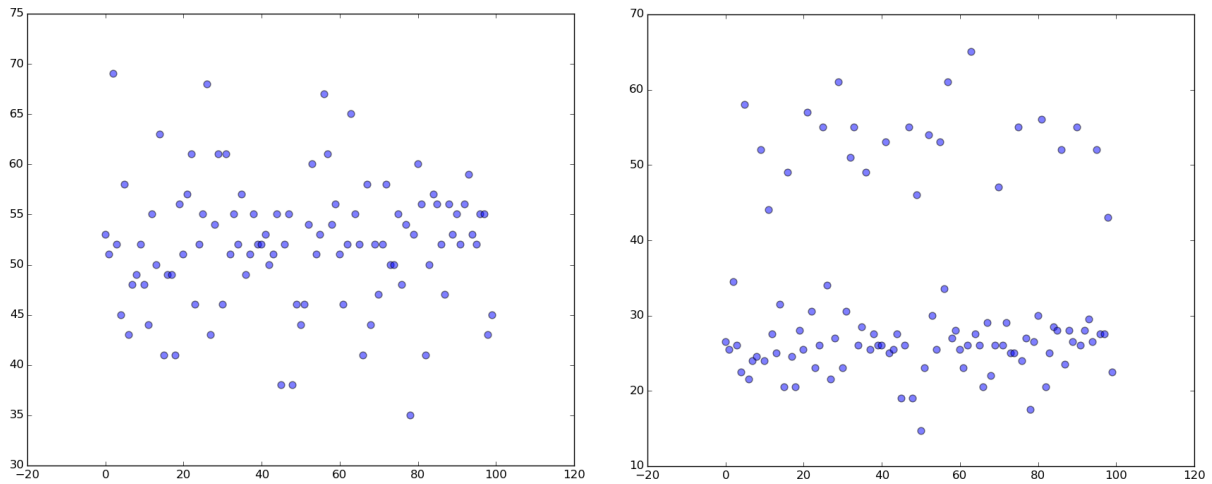


Figure 16: Trained model. Left - rewards for 1000 iterations; Right - average rewards for 1 finished game.

This graphs shows different behavior. Mean rewards are around 32 which is higher than the random action model results.

In the the Table 1 below statistical information about the models is presented:

	Random rewards	Random averages	Trained rewards	Trained averages
Min	-6	-0.85	35	14.67
Max	4	1.0	69	65.0
Mean	-0.62	-0.08	52.0	32.32
Median	-1.0	-0.15	52.0	27.0
STD	2.35	0.42	6.32	12.42
Variance	5.53	0.18	40.0	154.19

Table 1: Comparing

Figures 17 and 18 below ⁵⁹ shows relation between results for both models on the same scale. The top figure within each plot is for trained model, and the bottom one is for random action model. The median is marked with a bold dotted line on each figure.

This diagrams shows that the results of the trained model are higher then results of the random action model.

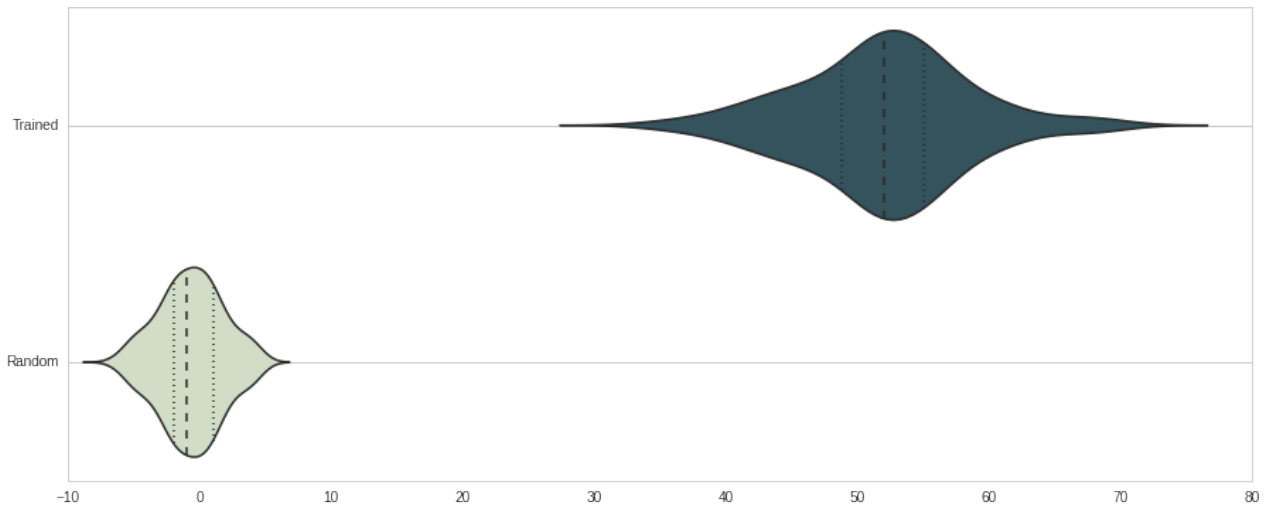


Figure 17: Compared rewards

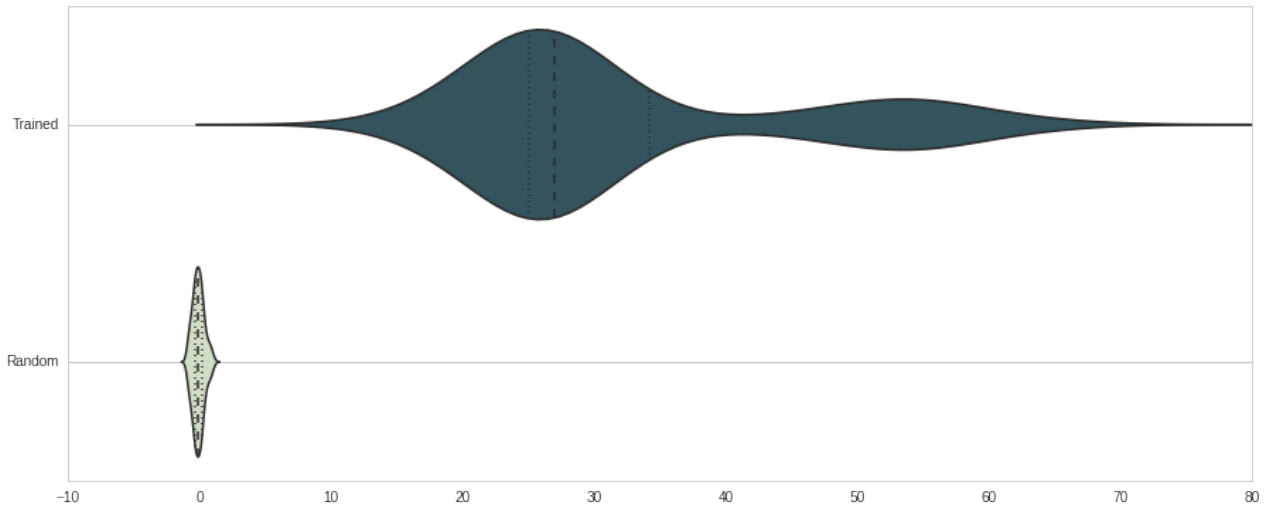
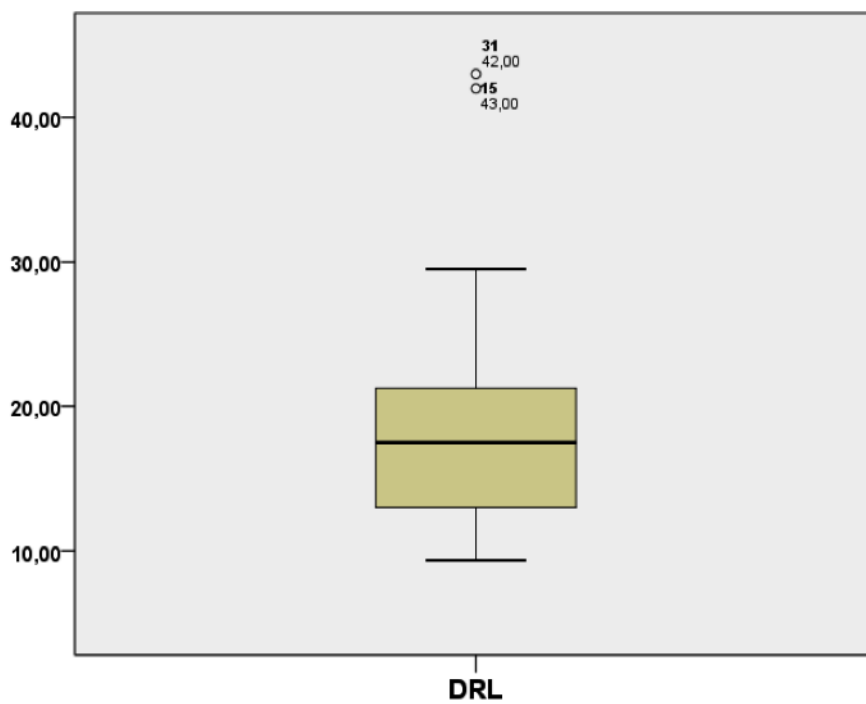
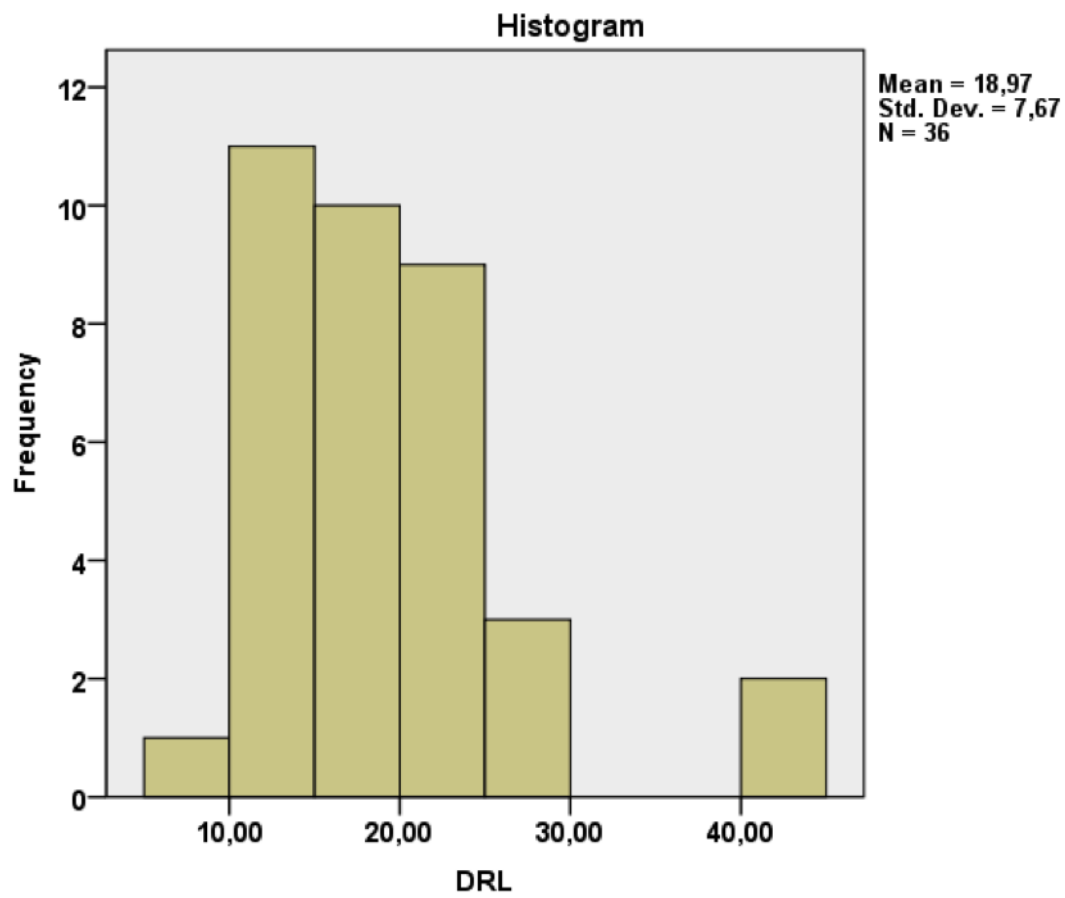


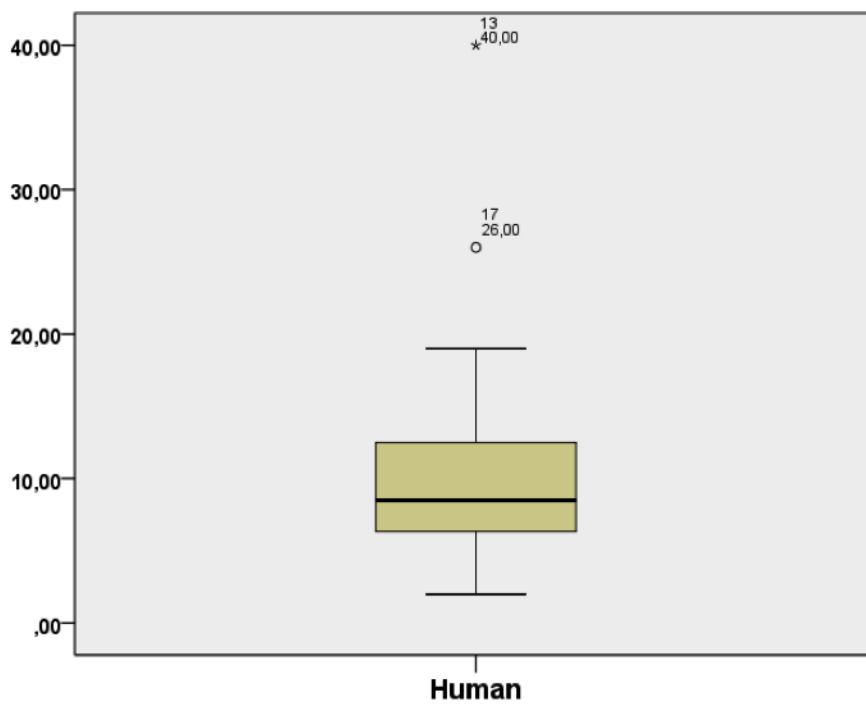
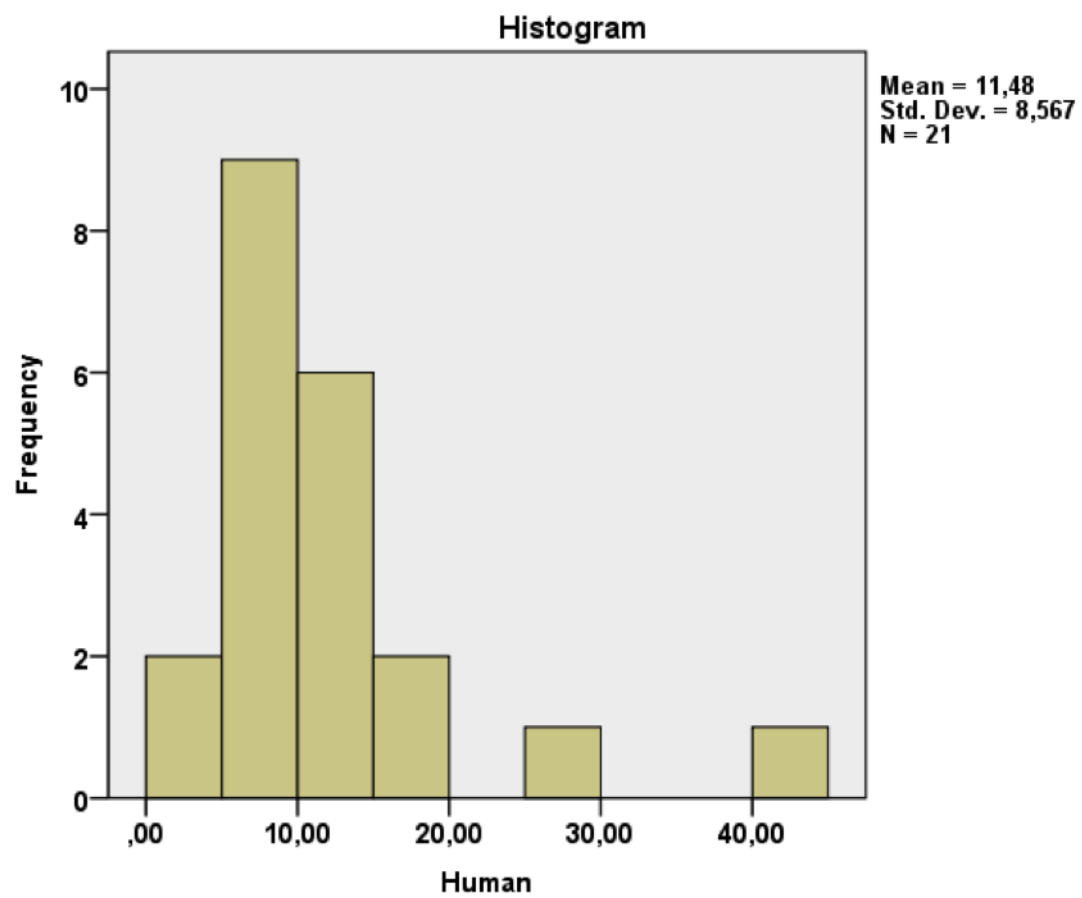
Figure 18: Compared averages

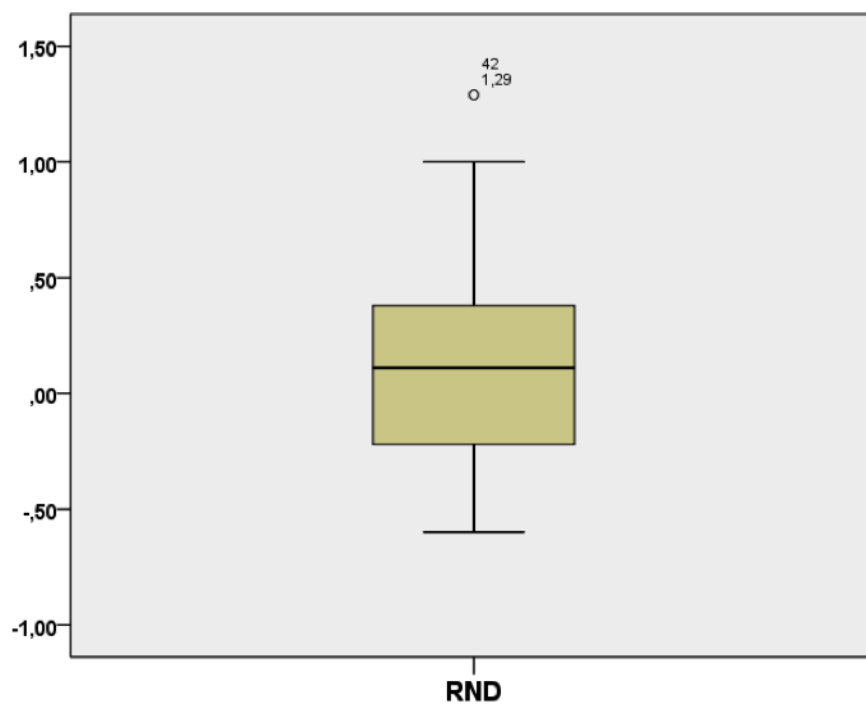
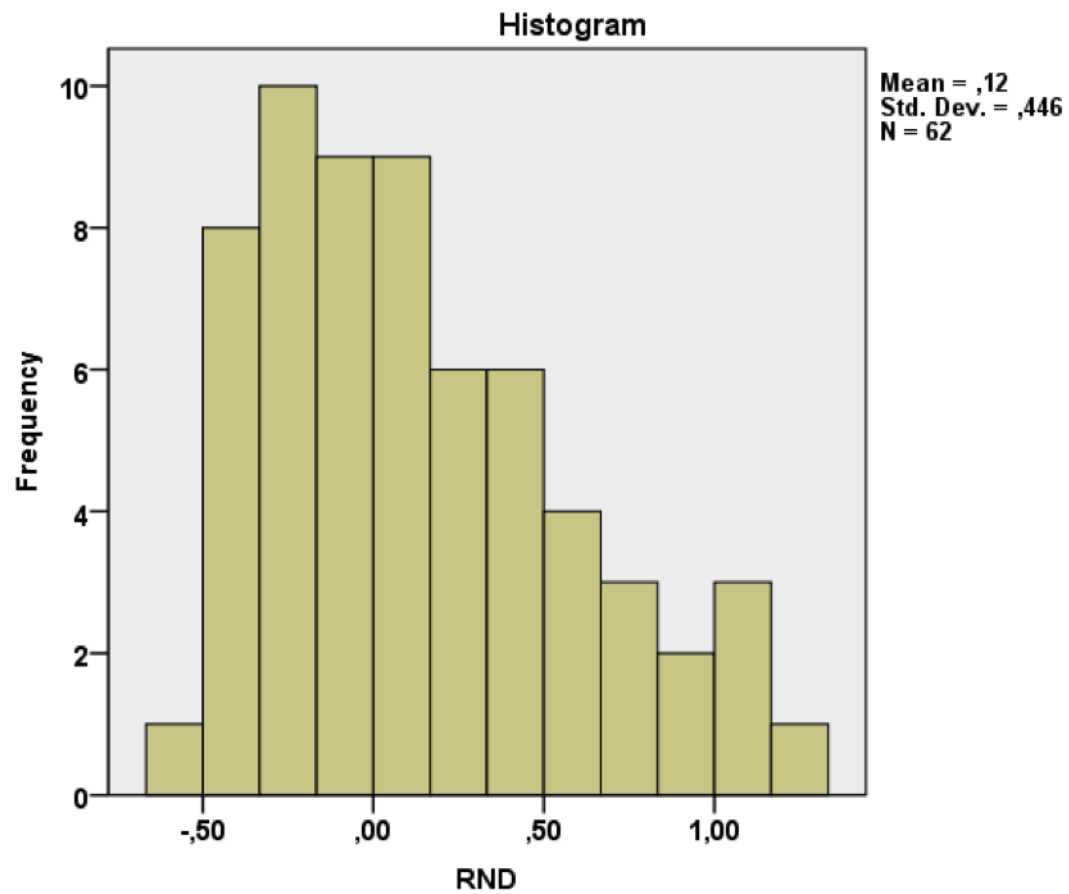
⁵⁹http://en.wikipedia.org/wiki/Violin_plot

7 Annex A1

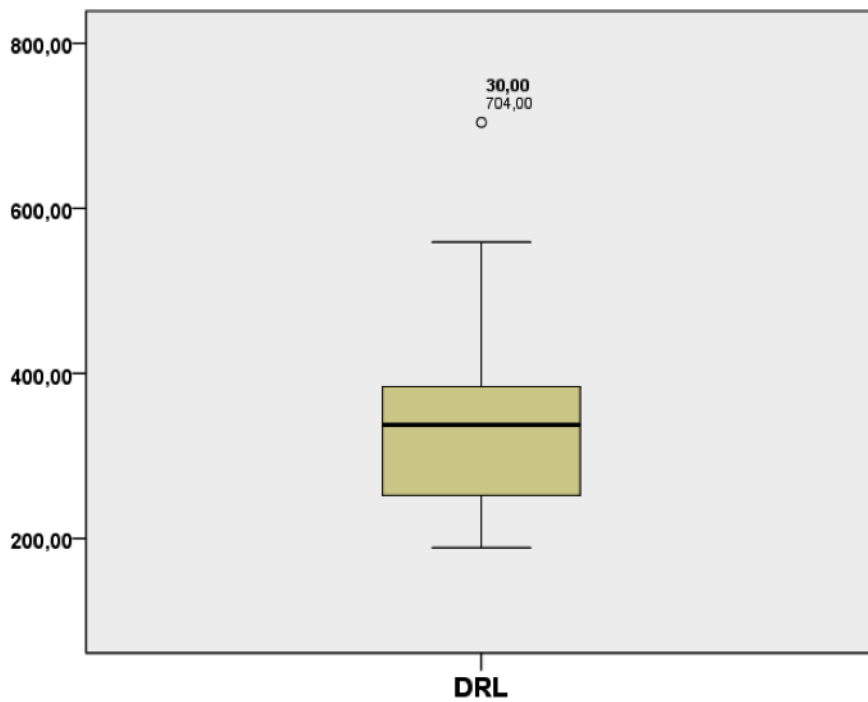
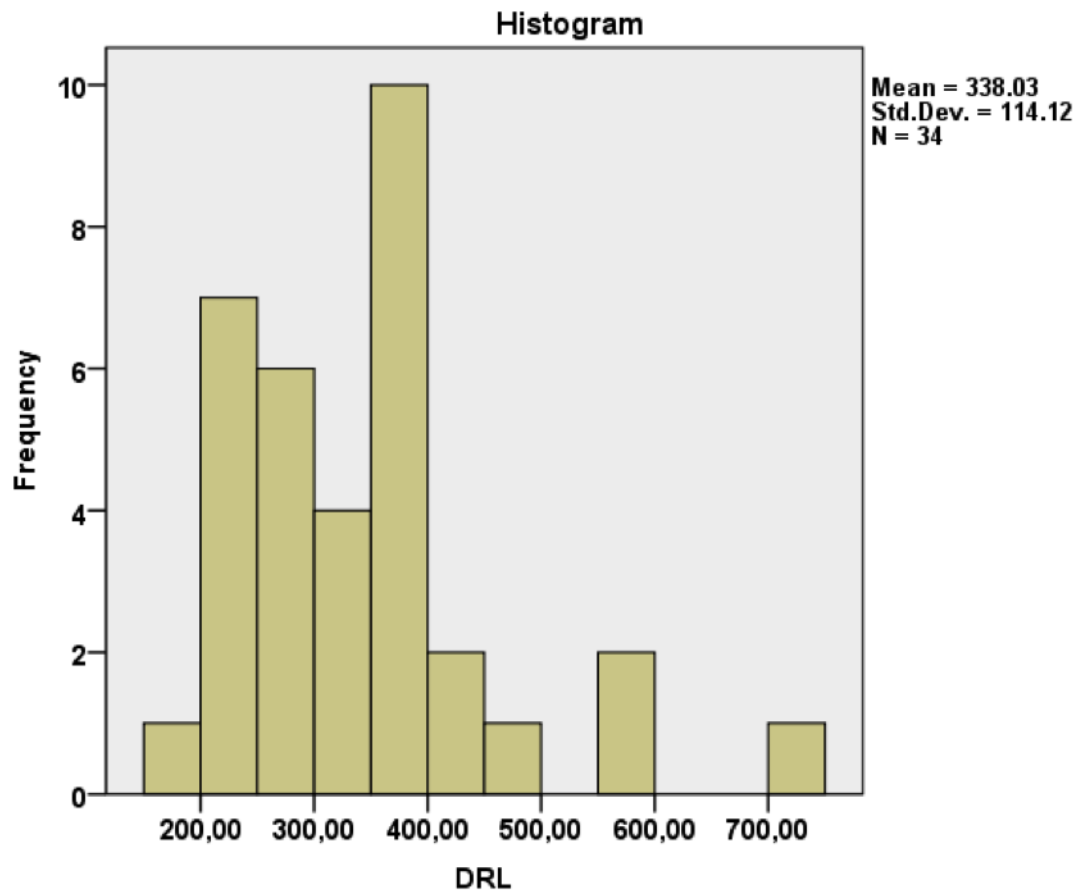
Breakout:

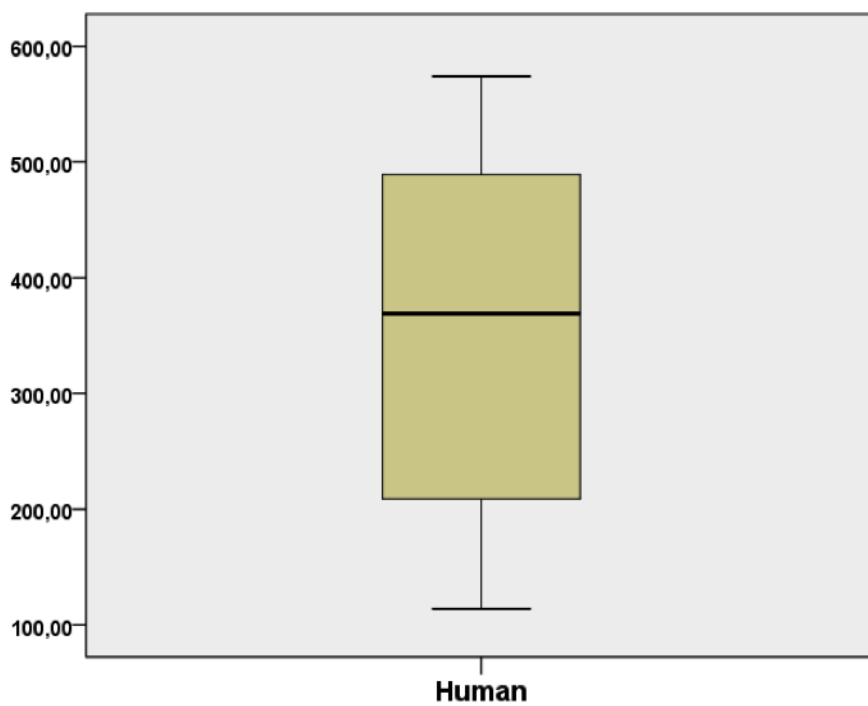
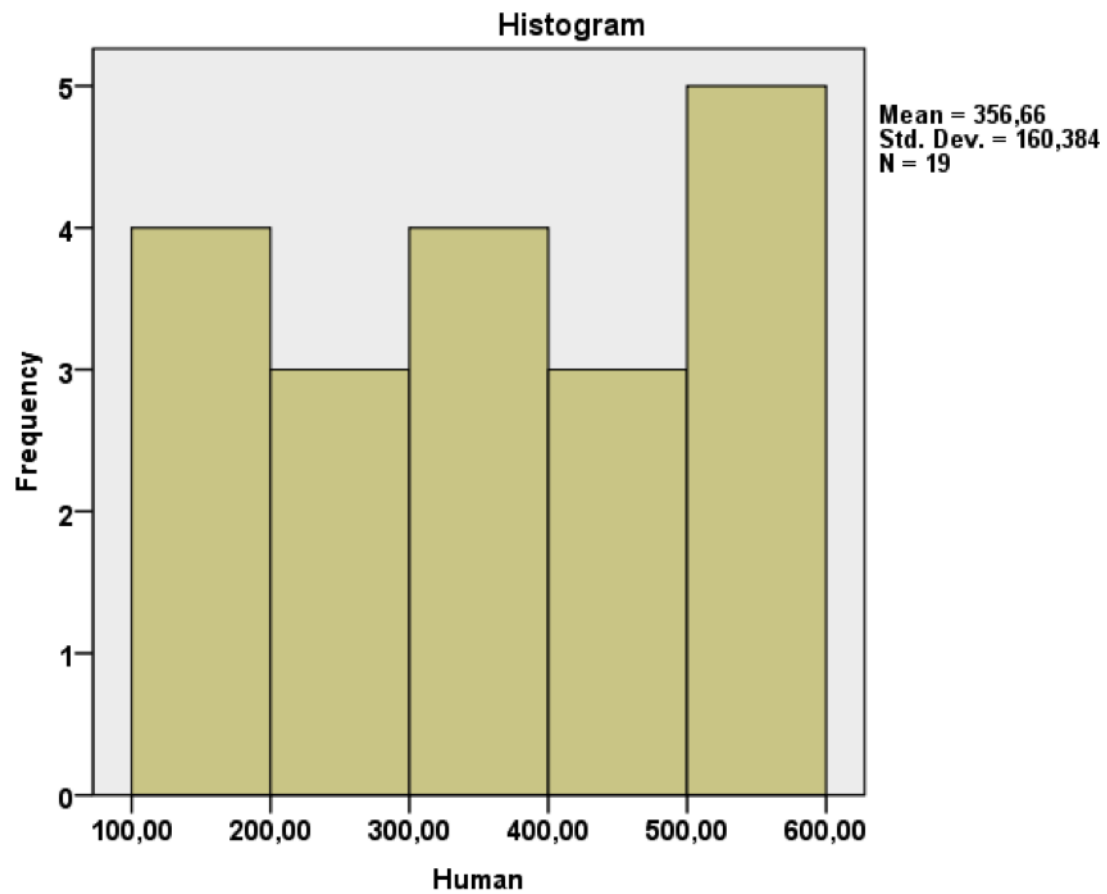


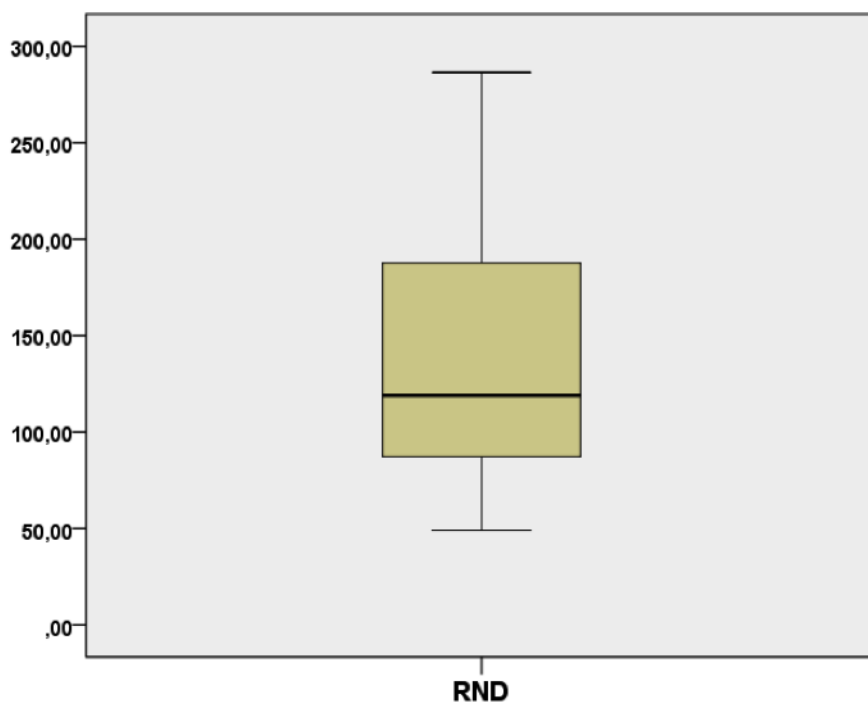
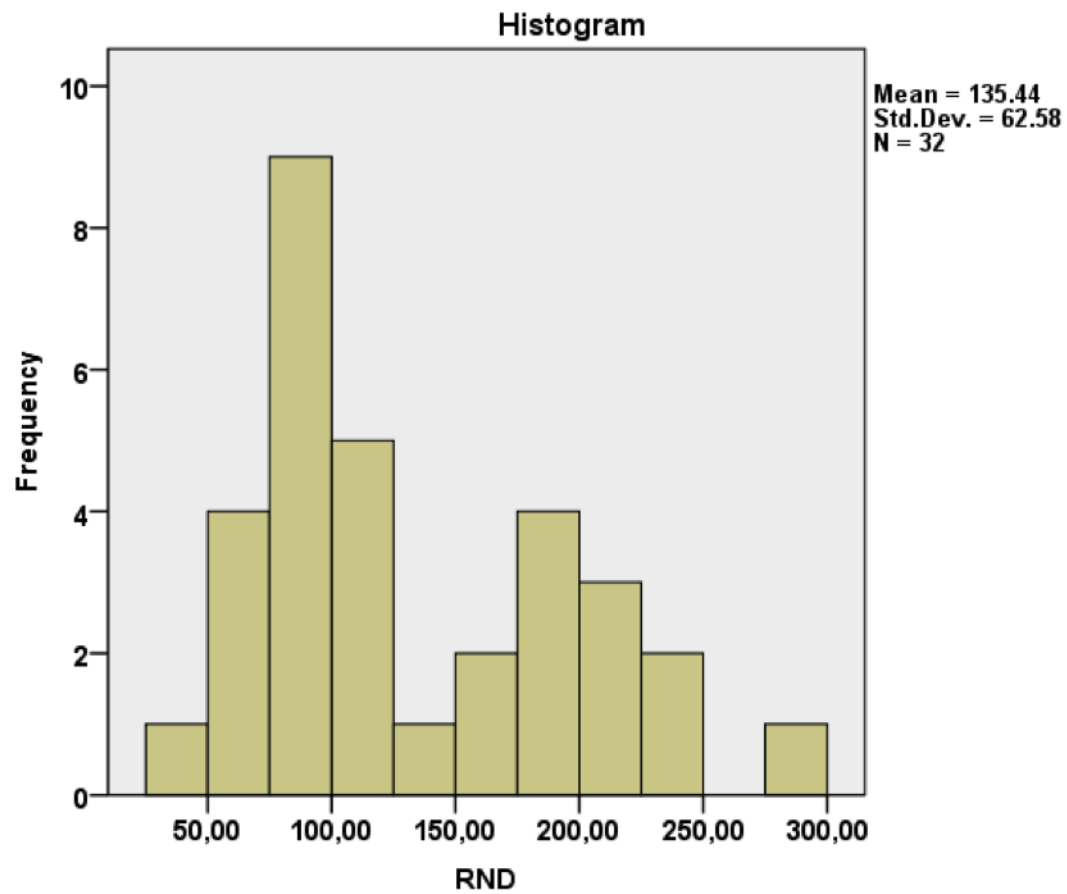




Space invaders:







Boxing:

