

# Raport Assignment 5

## Flappy Bird using Q learning

Membri echipei: Repede Monica-Gabriela, Murgulet Andrei

Acest raport are ca scop documentarea experimentelor pe care le-am facut pentru a realiza un agent inteligent pentru jocul Flappy Bird, cat mai performant, utilizand Q-learning.

### **Modelul 1 (11,1)**

#### **Arhitectura modelului**

Reteaua neuronală folosită pentru aproximarea valorilor Q are următoarea structură:  
input: imagini preprocesate  $66 \times 66 \times 4$  (4 imagini consecutive de  $66 \times 66$ )

-Stratul 1: convolutional cu 16 filtre de  $8 \times 8$  (kernel), stride de 4 și funcția de activare Relu

-Stratul 2: convolutional cu 32 filtre de  $4 \times 4$ , Stride 2 și funcția de activare Relu

-Stratul 3: Fully connected, out: 256 de neuroni și funcția de activare Relu

-Stratul 4: De ieșire, strat complet conectat liniar, cu o singură ieșire pentru fiecare acțiune (2 acțiuni, 0 și 1)

-> această arhitectură a rețelei este preluată din lucrarea "Playing Atari with Deep Reinforcement Learning" - [link document](#)

Preprocesarea imaginilor a constat din

-resize-ul acestora în (84, 110) urmat de un crop a imaginii obținute în dimensiunea (66,66)-> se elimină pământul, spațiul liber de dinaintea pasării și se limitează vederea următoarelor pipe-uri

-conversie în tonuri de gri (COLOR\_BGR2GRAY) pentru eficiență.

#### **Hyperparametri**

- episodes = 100 000 - numărul de episoade pe care se antrenează agentul
- experience\_memory = 50000 - dimensiunea pentru Experience Replay pentru a stoca cele mai recente frame-uri
- target\_update = 10000 numărul de pași între actualizările rețelei
- update\_freq = 4 - numărul de pași între actualizarea modelului Q-network
- frame\_skip = 3 - atunci când alegem să sarim(1), se da skip la 3 frame-uri (și se alege 0- adică nu face nimic)
- epsilon ia valori între 1 și 0.1 liniar timp de 1 milion de pași (alegeri de 0 sau 1) și apoi menținut la 0.1
- RMSProp cu mini batchuri de 32

**Sistemul de recompense:** pe lângă reward-ul din joc propriu-zis, urmărim unele scenarii:

1. pasarea se loveste in primul pipe(tinde sa faca salturi repetate pentru a nu se lovi de pamant si a primi o recompensa mai mare; la primul pipe ajunge la 101)->pentru a evita acest comportament de supravietuire neindicat, o penalizam cu -100
2. daca trece de 101 crestem cu +10

### **Rezultate obtinute:**

Pentru aceste date pe train in 100 000 de episoade rewardul maxim este de 1519 ceea ce in pipe-uri trecute reprezinta 38.

Pe test obtinem rezultate mai facile maximul fiind de 64 de pipe-uri trecute deoarece eliminam elementul de alegere random care persista in continuare pe train(prob 10%). Media insa este in jur de 25 de pipe uri (la 20 de episoade)

### **Modelul 2 (3,1)**

Acelasi model, insa **hyperparametrii** au suferit usoare modificari:

target\_update= 1000

politica de actualizare a epsilonului: pentru inceput am pornit cu modificare liniara pentru 100 000 de pasi, apoi cand a ajuns la episodul 38 000, pasul 516613, l-am modificat pentru 1 milion de pasi (epsilon vechi=0.1 , epsilon nou=0.54)

### **Rezultate obtinute:**

Pentru aceste date pe train in 100 000 de episoade rewardul maxim este de 1690 ceea ce in pipe-uri trecute reprezinta 43.

Pe test obtinem rezultate mai facile maximul fiind de 113 pipe-uri trecute. Media insa este tot in jur de 30 de pipe uri (la 20 de episoade)

*! De mentionat este ca in ambele cazuri rulara a mai fost oprita si astfel s-a pierdut memoria din bufferul de experience replay, retinandu-se in memorie doar modelul DQN pentru cand s a reluat rulara, iar de aici sunt posibile discrepante in continuitatea scorului !*

### **Modelul 3 (12,1)**

#### **Arhitectura modelului**

Reteaua neuronală folosită pentru aproximarea valorilor Q are următoarea structură:

-> [link paper](#)

input: imagini preprocesate 110\*84\*4 (4 imagini consecutive de 110\*84)

-Stratul 1: convolutional cu 32 filtre de 8x8(kernel), stride de 4 și funcția de activare Relu

-Stratul 2: convolutional cu 64 filtre de 4x4, Stride 2 și funcția de activare Relu

-Stratul 3: convolutional cu 64 filtre de 3x3, Stride 1 și funcția de activare Relu

-Stratul 4 ascuns: Fully connected, out: 512 de neuroni și funcția de activare Relu

-Stratul 5: De ieșire, strat complet conectat liniar, cu o singură ieșire pentru fiecare acțiune (2 acțiuni, 0 și 1)

Preprocesarea imaginilor a constat din

1. crop-ul imaginii - se elimină pământul, spațiul liber de dinaintea pasării, iar în imagine apare doar un pipe; al 2-lea pipe apare abia când pasarea este într-un pipe, aproape de final

2. resize-ul acestora în (84, 110)

3. conversie în tonuri de gri(COLOR\_BGR2GRAY) pentru eficiență

#### **Hyperparametri**

- episodes = 11 000 - numărul de episoade pe care se antrenează agentul
- experience\_memory = 50000 - dimensiunea pentru Experience Replay pentru a stoca cele mai recente frame-uri
- target\_update = 10000 numărul de pași între actualizările rețelei
- update\_freq = 4 - numărul de pași între actualizarea modelului Q-network
- frame\_skip = 3 - atunci când alegem să sarim(1), se da skip la 3 frame-uri (și se alege 0- adică nu face nimic)
- epsilon ia valori între 1 și 0.0001 liniar timp de 1 milion de pași(alegeri de 0 sau 1) și apoi menținut la 0.0001
- RMSProp cu mini batchuri de 64

**Sistemul de recompense:** pe lângă reward-ul din joc propriu-zis, urmărim unele scenarii:

1. pasarea se lovește în primul pipe(tinde să facă salturi repetate pentru a nu se lovi de pământ și a primi o recompensă mai mare; la primul pipe ajunge la 101)->pentru a evita acest comportament de supraviețuire neindicat, o penalizăm cu -50

2. dacă trece de 101, însă moare, penalizăm cu -10

3. dacă nu a murit și a trecut de 101, adunăm +1

4. dacă nu a murit însă nu a trecut de primul pipe (101) adunăm +0.1

## Rezultate obtinute:

Pentru aceste date pe train in 11 000 de episoade rewardul maxim este de 8720 ceea ce in pipe-uri trecute reprezinta 233.

Pe test obtinem rezultate mai facile maximul fiind de 387 pipe-uri trecute.

De asemenea am rulat de 6 ori cate 20 de episoade pentru a observa performanta agentului si am trecut scorurile (in functie de numarul de pipe-uri) in tabelul de mai jos.

	A	B	C	D	E	F	G	H	I							
	Tabelul_1															
1	#	Nr_ep	#	Incarcare 1	#	Incarcare 2	#	Incarcare 3	#	Incarcare 4	#	Incarcare 5	#	Incarcare 6		
2		1		44		104		49		47		23		43		
3		2		71		59		21		29		3		23		
4		3		34		117		6		31		201		39		
5		4		94		10		141		1		4		12		
6		5		53		3		9		72		9		11		
7		6		77		79		59		40		37		9		
8		7		54		3		44		69		67		3		
9		8		27		31		17		8		3		97		Average
10		9		25		141		111		29		91		121		49,55
11		10		73		25		149		93		5		109		
12		11		19		111		89		44		27		49		
13		12		16		28		77		2		57		56		
14		13		19		19		6		49		42		55		
15		14		4		18		105		100		89		9		
16		15		27		161		2		27		26		13		
17		16		3		4		37		237		149		40		
18		17		44		17		20		2		19		171		
19		18		116		9		107		11		48		20		
20		19		5		13		70		13		76		19		
21		20		23		129		19		7		24		89		
22				41,4		54,05		56,9		45,55		50		49,4		
23																

Dupa cum se poate observa, modelul din urma a produs cele mai bune rezultate atat pe train cat si pe test; acesta este modelul la care am ramas.