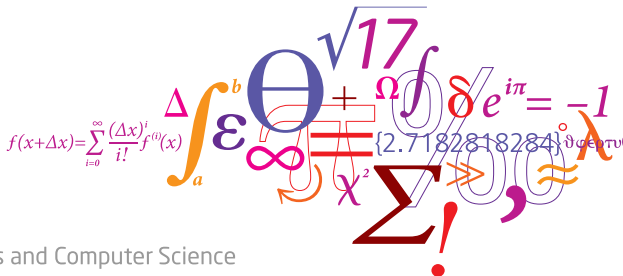


Deep Q learning

Ellen Vanhove, Emil Tyge, Carl Kjærgaard



DTU Compute

Department of Applied Mathematics and Computer Science

Outline

- Introduction
 - Our Problem
- Q-Learning
 - Reinforcement Learning
 - Different R-Learning Methods
 - Applying Deep Learning Methods
- Our Implementation
 - Breakout
 - Modelling
 - Difficulties
 - Plots
 - Comparison to other solutions

Why study these problems?

- Solutions to complex problems - State explosions
- Model-free AI planners
- Other real-value solutions

Problem Modelling

- Markovian Decision Problems
- On- vs Off-policy learning
 - ① Policies in Reinforcement Learning
 - ② Problems with on-policy learning

On Policy methods

- Policy Gradients

Off Policy Methods

- Q-Learning
- n-step Q-Learning

Mixed Policy Methods

- Actor-Critic methods
- A3C

Q Function Estimators

- Estimating complex linear functions
- Using deep networks as q-function estimator

Q-Learning

- Use ε -greedy policy

$$\pi(a|s, \varepsilon) = \begin{cases} \operatorname{argmax}_a Q(s, a, \theta) & \varepsilon > rand \\ a_{rand} & otherwise \end{cases}$$

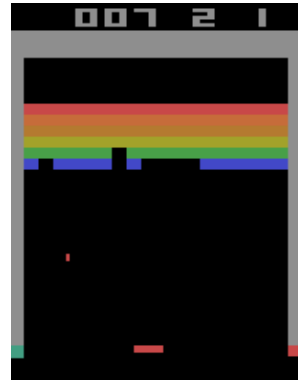
- Store s_n, a_n, r_n in replay memory
- Optimize

$$R_n = r_n + \gamma Q(s_{n+1}, \pi(s_{n+1}, 0), \theta^*)$$

$$\min_{\theta} \sum_n (R_n - Q(s_n, a_n, \theta))^2 + \alpha (R_n - Q(s_n, a_n, \theta))$$

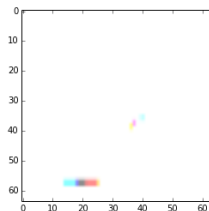
Breakout

- Hidden markov model - fully observable with multiple frames
- Large observation space
- 4 Actions
- Reward when scoring



Generalized input

- Generalized input
64x64 images, 3 temporal channels
- Convolutional Neural net
 - Providing location invariance and pattern matching
 - Reduces parameters to be learned
 - Learning element relationships
- Dense neural net
 - Collecting filter information
- Quadratic cost - with L2 regularization



Problems

- Complicated policies to learn, very general problem description
- Robustness and convergence of solutions
- Connecting actions to rewards
- Hyper parameter - long training time
- Large observation space - smaller state space

Our solutions

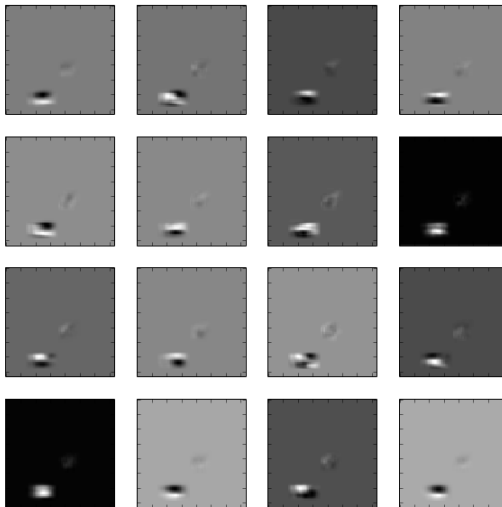
- Large memory of state/action/reward samples at least 250k
- Lookahead - calculating the discounted rewards
- Appropriate rewards - punish loosing the ball as well
- Cropping - simplifying the observational space
- L2 parameter costs - restrain parameters

Evolution of Q values when training



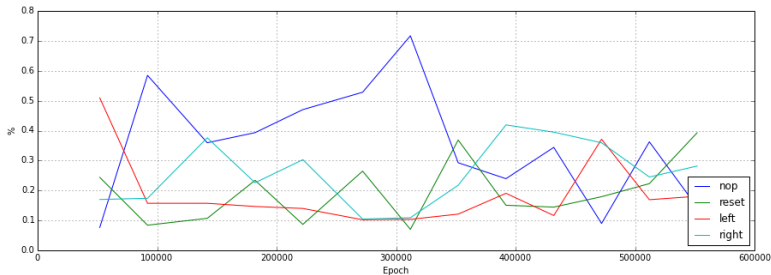
Our Implementation

Filtered frame



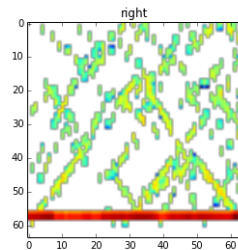
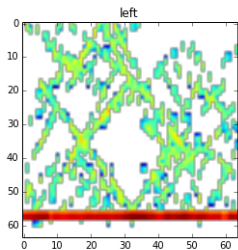
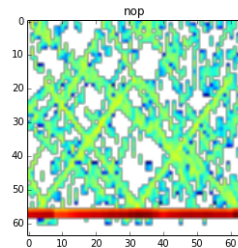
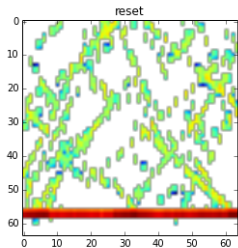
Our Implementation

Performance testing



Our Implementation

Analysing policies



Solutions

- Deepmind
- A3C