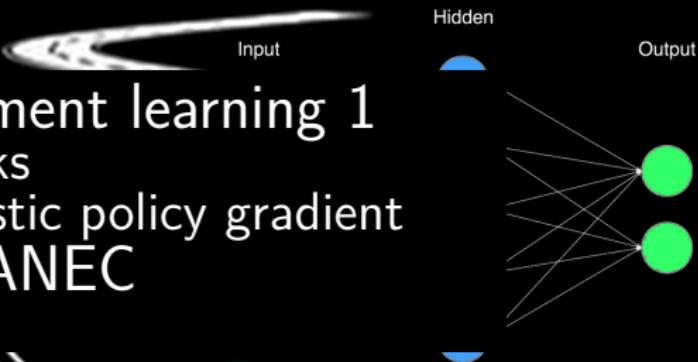


$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)]$$

(The New Action Value = The Old Value) + The Learning Rate \times (The New Information — the Old Information)



deep reinforcement learning 1

- deep Q networks
- deep deterministic policy gradient

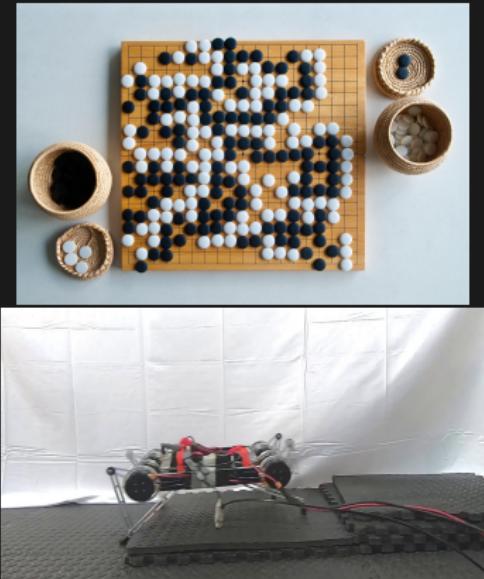
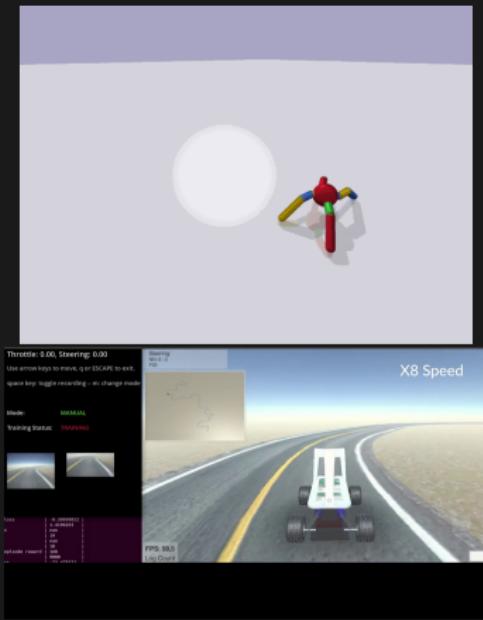
Michal CHOVANEC



reinforcement learning

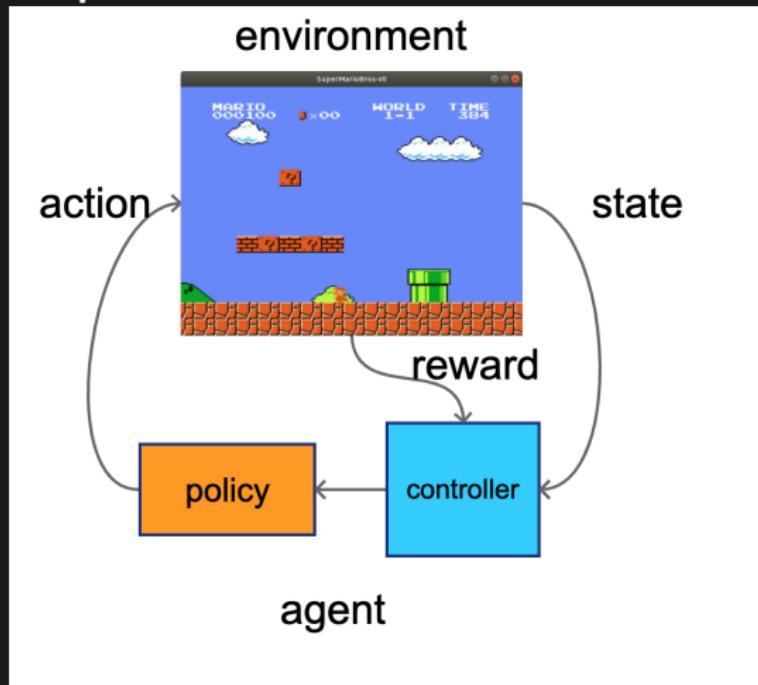


reinforcement learning

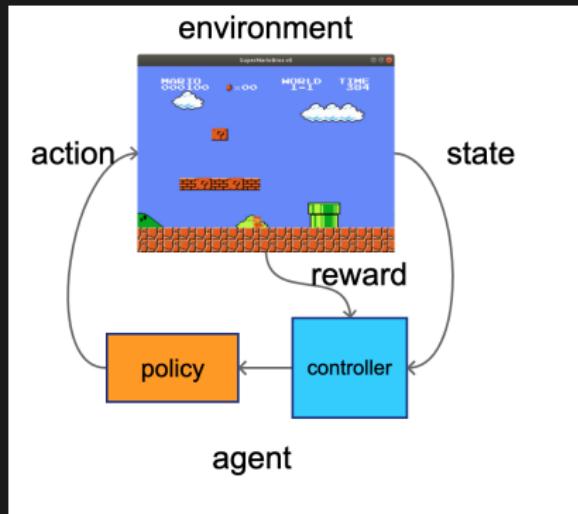


reinforcement learning

learning from punishments and rewards



reinforcement learning



- obtain state
- select action
- execute action
- learn from experiences

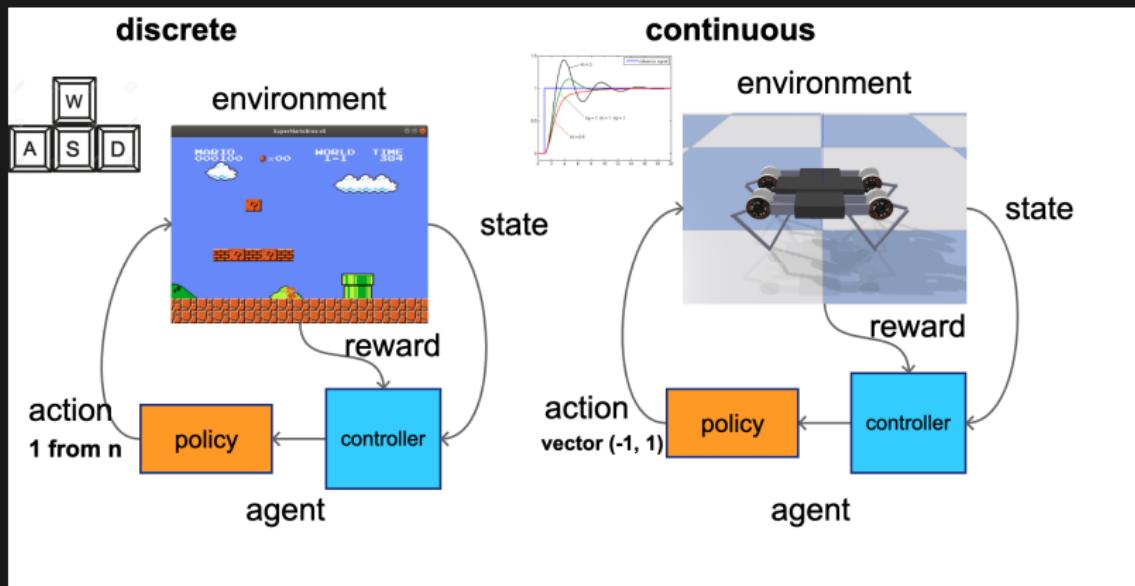
reinforcement learning - algorithms

- discrete actions space
 - Deep Q-network, DQN
 - Dueling DQN
 - Rainbow DQN
- continuous actions space
 - Actor Critic
 - Advantage Actor Critic
 - Proximal policy optimization
 - Soft Actor critic
 - Deep deterministic policy gradient
 - D4PG, SDDPG
- model based
 - curiosity
 - world models
 - imagination agents

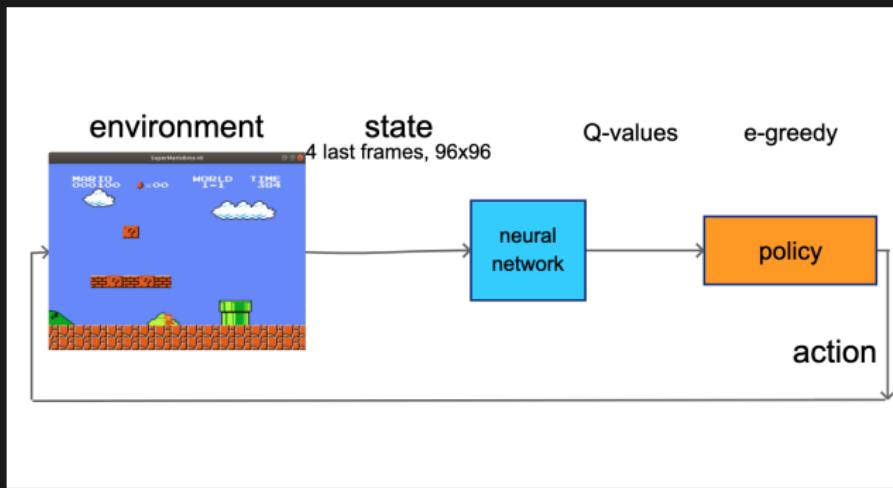
f.e. SDDPG - sampled DDPG, based on Wasserstein loss : Optimal transport, Cédric Villani, 600+ pages

action space

- discrete action space
 - keys, keypad
- continuous action space
 - motors, PWMs, steering, force control



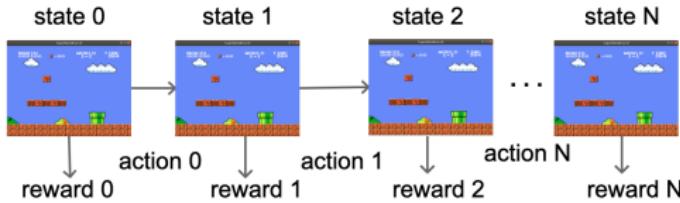
deep Q learning



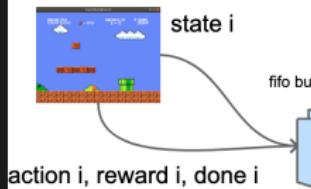
- ① play games
- ② store transitions into buffer
 - state, action, reward, done
- ③ learn from buffer

deep Q learning

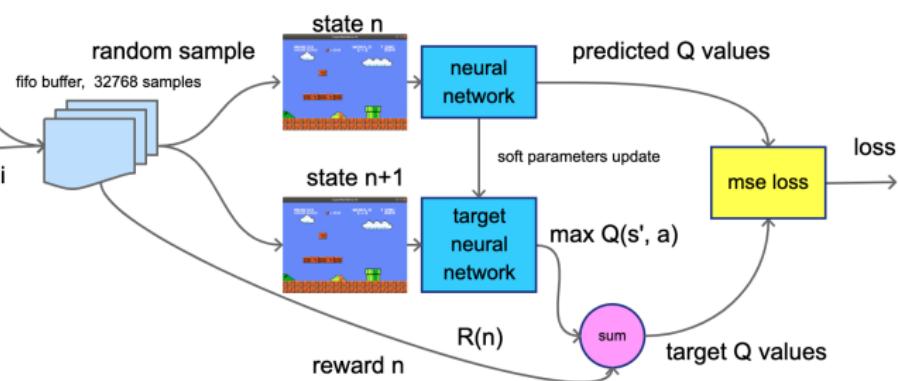
1, game play



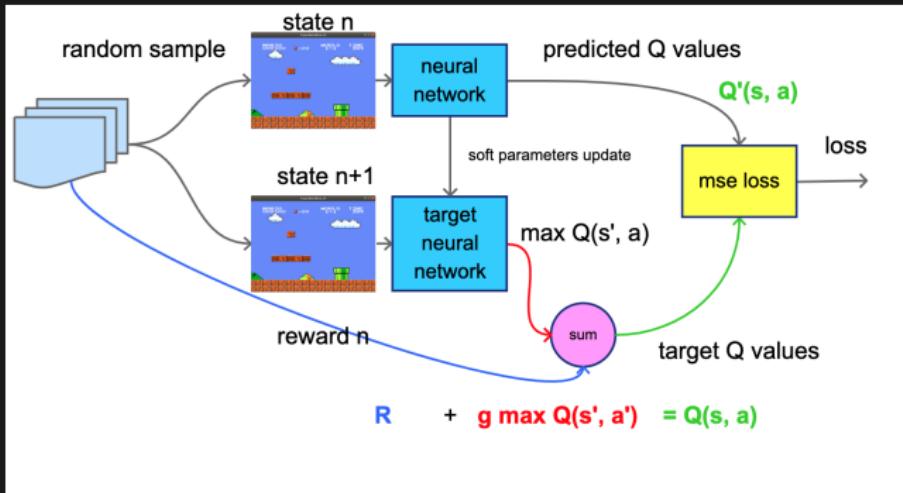
2, experience replay buffer



3, train network



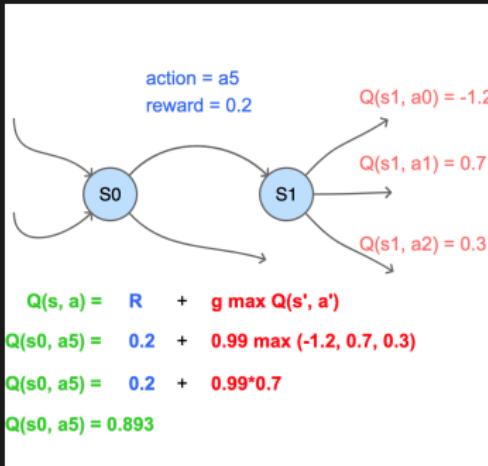
deep Q learning



$$Q(s, a; \theta) = \underset{\text{reward}}{R} + \gamma \max_{a'} \underset{\text{discounted future reward}}{Q(s', a'; \theta^-)}$$

$$\mathcal{L}(\theta) = \left(R + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2$$

deep Q learning



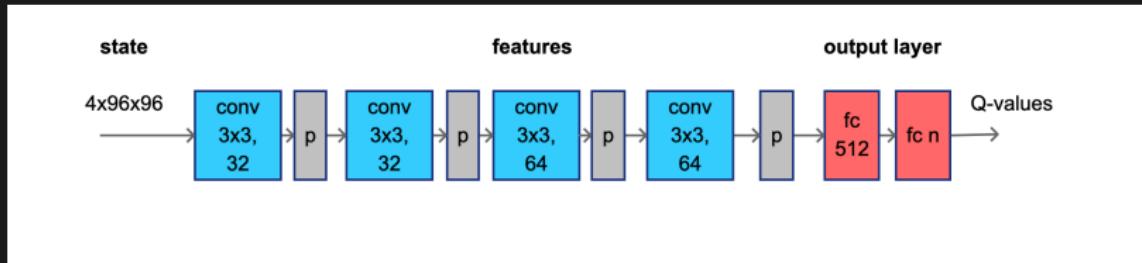
$$Q(s, a; \theta) = \underbrace{R}_{\text{reward}} + \gamma \max_{a'} Q(s', a'; \theta^-) \quad \text{discounted future reward}$$

$$\mathcal{L}(\theta) = \left(R + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta) \right)^2$$

dqn code

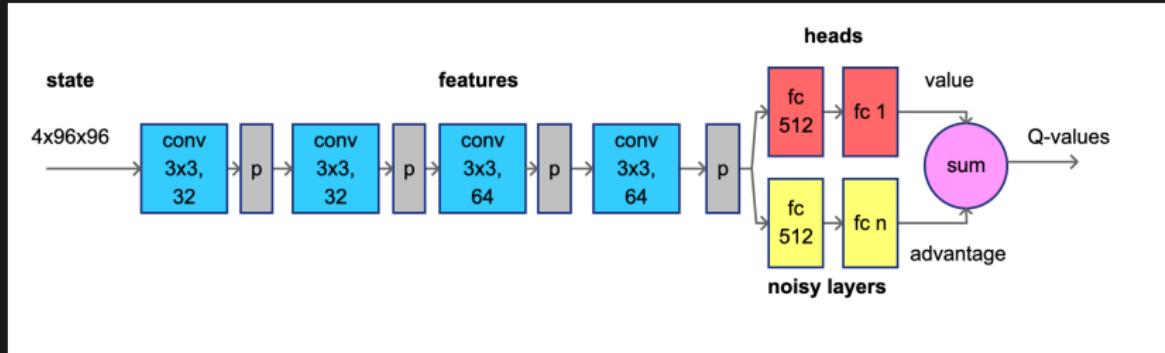
```
213     def train_model(self):
214         state_t, action_t, reward_t, state_next_t, done_t = self.experience_replay.sample(self.batch_size,
215
216         #q values, state now, state next
217         q_predicted      = self.model.forward(state_t)
218         q_predicted_next = self.model_target.forward(state_next_t)
219
220         #compute target, n-step Q-learning
221         q_target        = q_predicted.clone()
222         for j in range(self.batch_size):
223             action_idx          = action_t[j]
224             q_target[j][action_idx] = reward_t[j] + torch.max[q_predicted_next[j]*(1 - done_t[j])]
225
226
227         #train DQN model
228         loss = ((q_target.detach() - q_predicted)**2)
229         loss = loss.mean()
230
231         self.optimizer.zero_grad()
232         loss.backward()
233         for param in self.model.parameters():
234             param.grad.data.clamp_(-10.0, 10.0)
235         self.optimizer.step()
```

model architecture



- input 96x96 grayscale, 4 stacked frames
- 3x3 convs + pooling
- two fully connected layers
- small learning rate $\eta = 0.0001$, batch size = 32
- $\gamma = 0.99$
- exploration ϵ -greedy, 1M samples linear decay from 1 to 0.05
- total training 10M samples

dueling DQN, model architecture



$$Q(s, a) = V(s) + A(s, a)$$

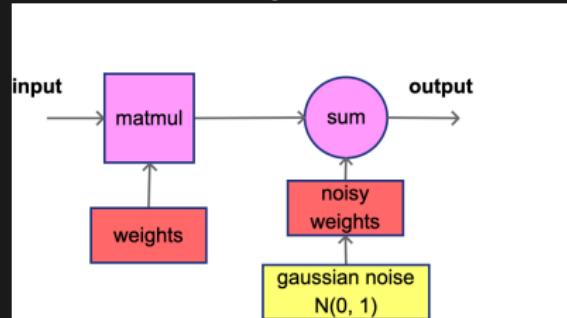
$$Q(s, a) = V(s) + A(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a' \in \mathcal{A}} A(s, a')$$

WRONG : `q = value + advantage - advantage.mean()`

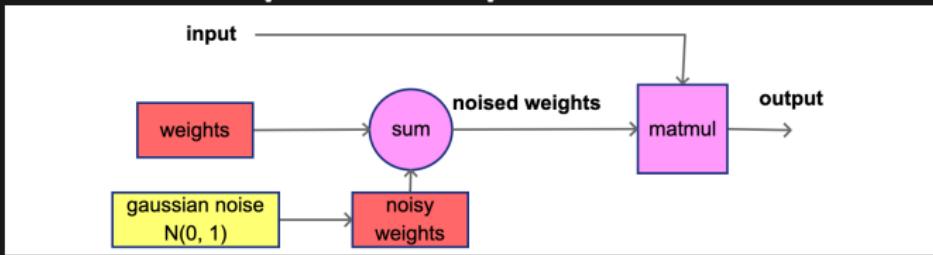
CORRECT : `q = value + advantage - advantage.mean(dim=1, keepdim=True)`

noisy layers for exploration

action space noise



parameter space noise



model architecture

```
self.layers_features = [
    nn.Conv2d(input_channels, 32, kernel_size=3, stride=1, padding=1),
    nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2, padding=0),

    nn.Conv2d(32, 32, kernel_size=3, stride=1, padding=1),
    nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2, padding=0),

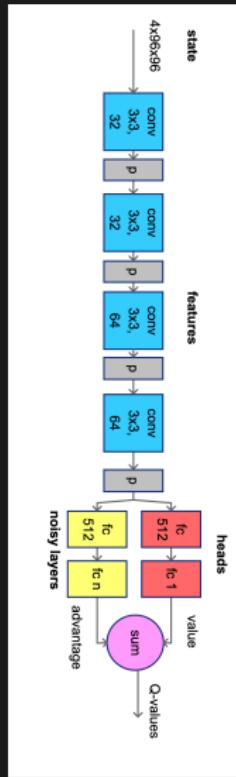
    nn.Conv2d(32, 64, kernel_size=3, stride=1, padding=1),
    nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2, padding=0),

    nn.Conv2d(64, 64, kernel_size=3, stride=1, padding=1),
    nn.ReLU(),
    nn.MaxPool2d(kernel_size=2, stride=2, padding=0),

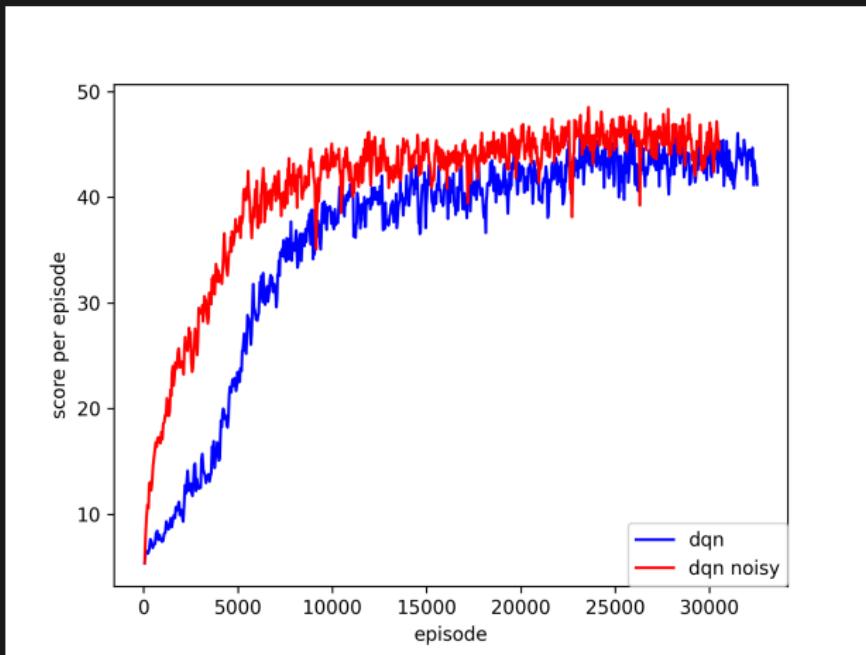
    Flatten()
]

self.layers_value = [
    nn.Linear(fc_inputs_count, 512),
    nn.ReLU(),
    nn.Linear(512, 1)
]

self.layers_advantage = [
    libs_layers.NoisyLinear(fc_inputs_count, 512),
    nn.ReLU(),
    libs_layers.NoisyLinear(512, outputs_count)
]
```



results on MsPacman

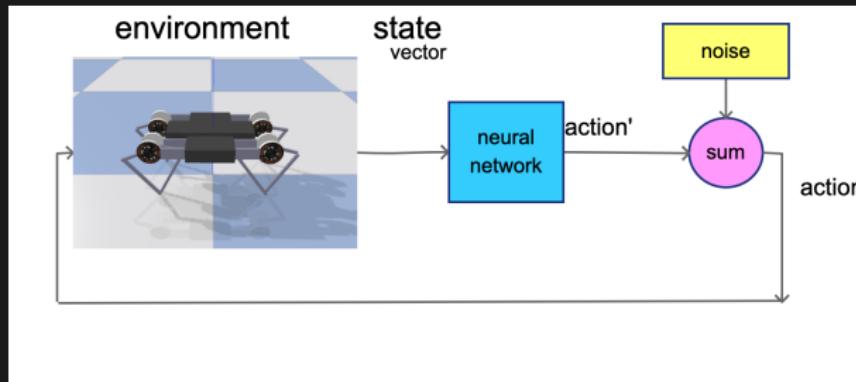


wise Wizard's DQN spell chart

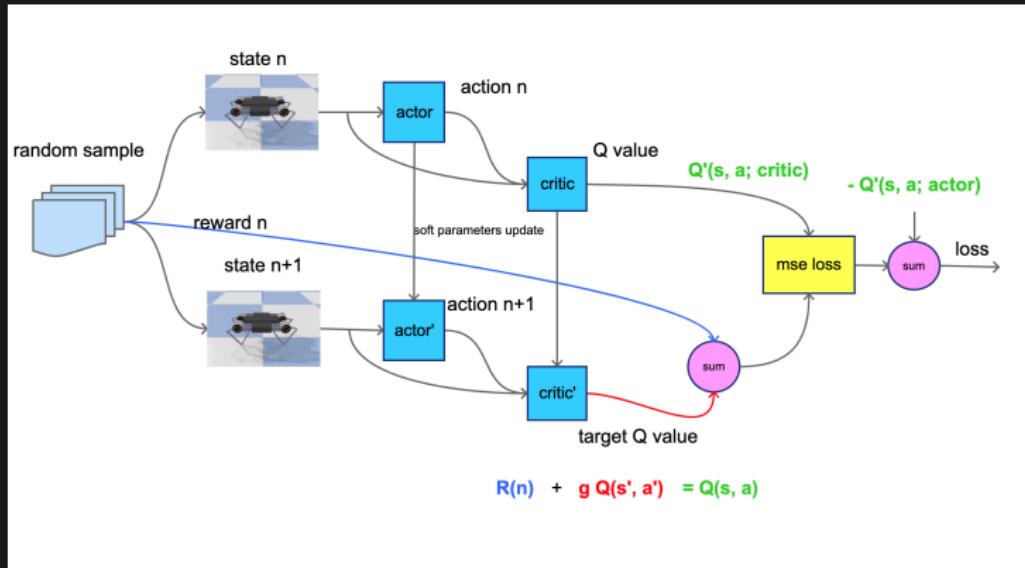
- use **input** with nice size $k2^j$, 64, 96, 128 ...
- **conv layers** use 3x3 convs, poolings or strides > 1
- for **weight init** use Xavier
- **don't use** batchnorm (use SkipInit - not tested yet),
- **don't use** weight regularisation
- use **soft** target network update, $\tau = 0.001$
- slow learning rate $\eta = 0.0001$ works better
- for exploration use ϵ -greedy decay from 1.0 to 0.05 in 1M steps, or use noisy layers

deep deterministic policy gradient

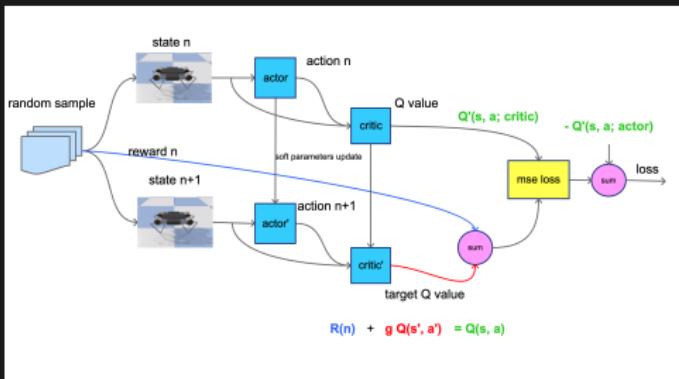
- continuous action space
- natural extension of DQN
- actor-critic structure



DDPG



DDPG



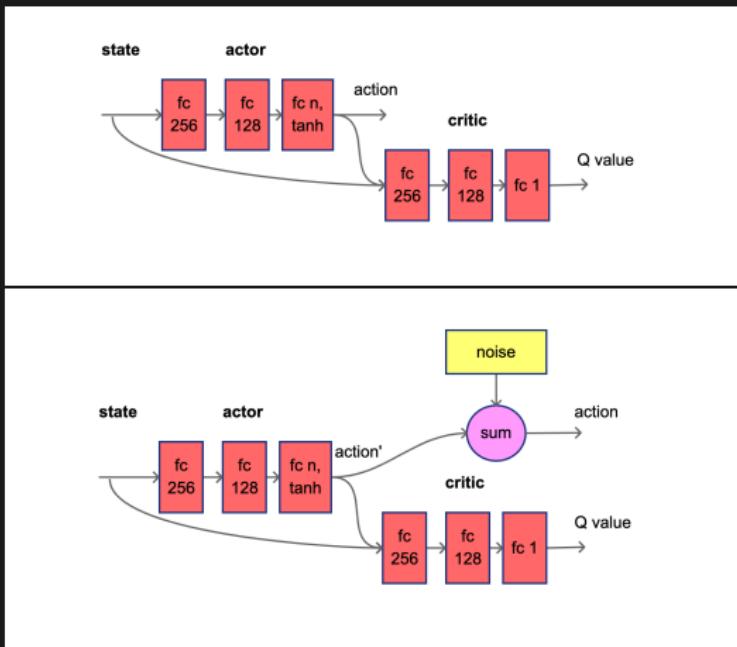
$$\mathcal{L}(\theta) = (R + \gamma Q(s', A(s'; \phi^-); \theta^-) - Q(s, A(s; \phi); \theta))^2$$

$$\mathcal{L}(\phi) = -Q(s, A(s; \phi); \theta)$$

where

- Q is critic network with parameters θ
- A is actor network with parameters ϕ

model architecture



ddpg actor code

```
class Model(torch.nn.Module):
    def __init__(self, input_shape, outputs_count, hidden_count = 256):
        super(Model, self).__init__()
        self.device = "cpu"

        self.layers = [
            nn.Linear(input_shape[0], hidden_count),
            nn.ReLU(),
            nn.Linear(hidden_count, hidden_count//2),
            nn.ReLU(),
            nn.Linear(hidden_count//2, outputs_count),
            nn.Tanh()
        ]

        torch.nn.init.xavier_uniform_(self.layers[0].weight)
        torch.nn.init.xavier_uniform_(self.layers[2].weight)
        torch.nn.init.uniform_(self.layers[4].weight, -0.3, 0.3)

        self.model = nn.Sequential(*self.layers)
        self.model.to(self.device)

        print(self.model)

    def forward(self, state):
        return self.model(state)
```

ddpg critic code

```
class Model(torch.nn.Module):
    def __init__(self, input_shape, outputs_count, hidden_count = 256):
        super(Model, self).__init__()

        self.device = "cpu"

        self.layers = [
            nn.Linear(input_shape[0] + outputs_count, hidden_count),
            nn.ReLU(),
            nn.Linear(hidden_count, hidden_count//2),
            nn.ReLU(),
            nn.Linear(hidden_count//2, 1)
        ]

        torch.nn.init.xavier_uniform_(self.layers[0].weight)
        torch.nn.init.xavier_uniform_(self.layers[2].weight)
        torch.nn.init.uniform_(self.layers[4].weight, -0.003, 0.003)

        self.model = nn.Sequential(*self.layers)
        self.model.to(self.device)

        print(self.model)

    def forward(self, state, action):
        x = torch.cat([state, action], dim = 1)
        return self.model(x)
```

ddpg code

```
83 def train_model(self):
84     state_t, action_t, reward_t, state_next_t, done_t = self.experience_replay.sample(self.batch_size, self.model_critic.device)
85
86     reward_t = reward_t.unsqueeze(-1)
87     done_t   = (1.0 - done_t).unsqueeze(-1)
88
89     action_next_t  = self.model_actor_target.forward(state_next_t).detach()
90     value_next_t   = self.model_critic_target.forward(state_next_t, action_next_t).detach()
91
92     #critic loss
93     value_target   = reward_t + self.gamma*done_t*value_next_t
94     value_predicted = self.model_critic.forward(state_t, action_t)
95
96     critic_loss    = ((value_target - value_predicted)**2)
97     critic_loss    = critic_loss.mean()
98
99     #update critic
100    self.optimizer_critic.zero_grad()
101    critic_loss.backward()
102    self.optimizer_critic.step()
103
104    #actor loss
105    actor_loss     = -self.model_critic.forward(state_t, self.model_actor.forward(state_t))
106    actor_loss     = actor_loss.mean()
107
108    #update actor
109    self.optimizer_actor.zero_grad()
110    actor_loss.backward()
111    self.optimizer_actor.step()
...  
...
```

wise Wizard's DDPG spell chart

- **neurons count** on 1st layer = 10x state vector size
- **neurons count** on 2nd layer = 0.5x neurons on 1st layer
- **weight init** for hidden layers : use Xavier
- **weight init** actor output : use uniform $\langle -0.3, 0.3 \rangle$
- **weight init** critic output : use uniform $\langle -0.003, 0.003 \rangle$
- **gaussian noise** : linear decay variance, from 0.5 to 0.1, for 1M steps, or noisy layers
- use **soft** target network update, $\tau = 0.001$
- actor learning rate $\eta_a = 0.0001$
- critic learning rate $\eta_c = 0.0002$

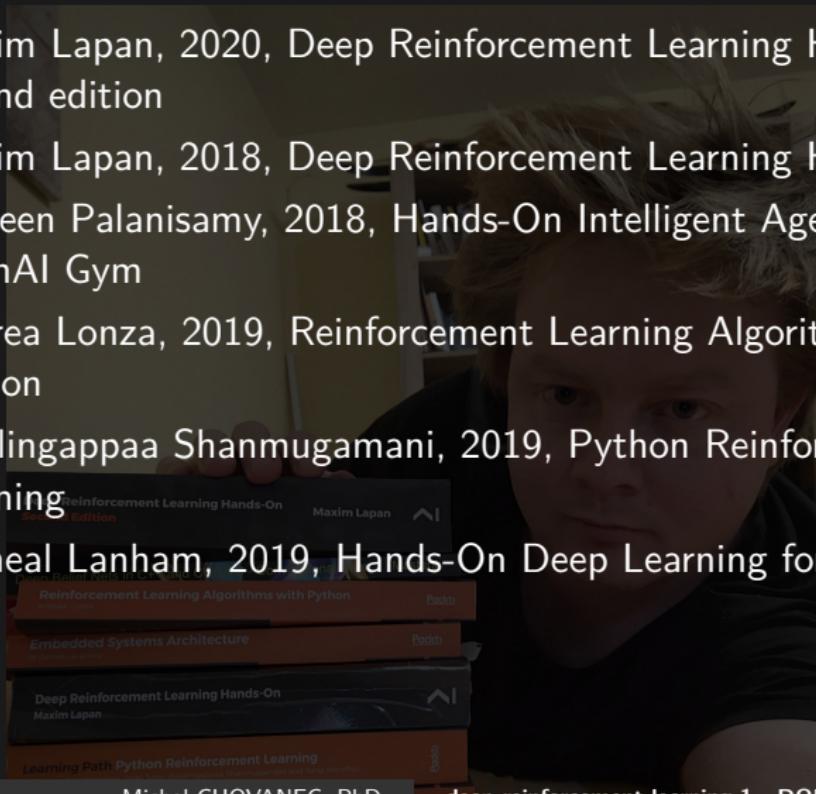
wise Wizard's magic staff

- fully connected nets (robotic envs) **train on CPU** - AMD Ryzen
- convolutional nets (visual inputs envs) **train on GPU**
- use fast CPU - envs are slow
- 32GB of RAM is enough
- for small visual envs (Atari, DOOM, Nec) - GTX1060, GTX1080ti, RTX2080 ...



books to read

- Maxim Lapan, 2020, Deep Reinforcement Learning Hands-On second edition
- Maxim Lapan, 2018, Deep Reinforcement Learning Hands-On
- Praveen Palanisamy, 2018, Hands-On Intelligent Agents with OpenAI Gym
- Andrea Lonza, 2019, Reinforcement Learning Algorithms with Python
- Rajalingappa Shanmugamani, 2019, Python Reinforcement Learning
- Micheal Lanham, 2019, Hands-On Deep Learning for Games



Q&A



michal.nand@gmail.com

https://github.com/michalnand/imagination_reinforcement_learning