NPFL122, Lecture 6



Rainbow

Milan Straka

i November 19, 2018





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated

Function Approximation

We will approximate value function v and/or state-value function q, choosing from a family of functions parametrized by a weight vector $w \in \mathbb{R}^d$.

We denote the approximations as

 $\hat{v}(s,oldsymbol{w}), \ \hat{q}(s,a,oldsymbol{w}).$

We utilize the Mean Squared Value Error objective, denoted \overline{VE} :

$$\overline{VE}(oldsymbol{w}) \stackrel{ ext{def}}{=} \sum_{s \in \mathcal{S}} \mu(s) \left[v_{\pi}(s) - \hat{v}(s,oldsymbol{w})
ight]^2,$$

where the state distribution $\mu(s)$ is usually on-policy distribution.

DDQN

NPFL122, Lecture 6

Refresh DQN

PriRep Duelling

NoisyNets I



Gradient and Semi-Gradient Methods



The functional approximation (i.e., the weight vector \boldsymbol{w}) is usually optimized using gradient methods, for example as

$$oldsymbol{w}_{t+1} \leftarrow oldsymbol{w}_t - rac{1}{2} lpha
abla \left[v_\pi(S_t) - \hat{v}(S_t,oldsymbol{w}_t)
ight]^2 \ \leftarrow oldsymbol{w}_t + lpha \left[v_\pi(S_t) - \hat{v}(S_t,oldsymbol{w}_t)
ight]
abla \hat{v}(S_t,oldsymbol{w}_t).$$

As usual, the $v_{\pi}(S_t)$ is estimated by a suitable sample. For example in Monte Carlo methods, we use episodic return G_t , and in temporal difference methods, we employ bootstrapping and use $R_{t+1} + \gamma \hat{v}(S_{t+1}, \boldsymbol{w})$.

DDQN

Deep Q Network



Off-policy Q-learning algorithm with a convolutional neural network function approximation of action-value function.

Training can be extremely brittle (and can even diverge as shown earlier).



Figure 1 of the paper "Human-level control through deep reinforcement learning" by Volodymyr Mnih et al.

NPFL122, Lecture 6

DQN DDQN

Refresh

PriRep

Duelling NoisyNets

DistRL

Rainbow

4/37

Deep Q Networks

- Preprocessing: 210×160 128-color images are converted to grayscale and then resized to 84×84 .
- Frame skipping technique is used, i.e., only every 4^{th} frame (out of 60 per second) is considered, and the selected action is repeated on the other frames.
- Input to the network are last 4 frames (considering only the frames kept by frame skipping), i.e., an image with 4 channels.
- The network is fairly standard, performing
 - $^{\circ}~$ 32 filters of size 8×8 with stride 4 and ReLU,
 - $^{\circ}~$ 64 filters of size 4×4 with stride 2 and ReLU,
 - $^{\circ}~$ 64 filters of size 3×3 with stride 1 and ReLU,
 - $^{\circ}\,$ fully connected layer with 512 units and ReLU,
 - $^{\circ}$ output layer with 18 output units (one for each action)

DDQN

NPFL122, Lecture 6

Duelling

Deep Q Networks

Ú FÁ

• Network is trained with RMSProp to minimize the following loss:

$$\mathcal{L} \stackrel{ ext{\tiny def}}{=} \mathbb{E}_{(s,a,r,s') \sim data} \left[(r + \gamma \max_{a'} Q(s',a';ar{ heta}) - Q(s,a; heta))^2
ight].$$

• An ε -greedy behavior policy is utilized.

Important improvements:

- experience replay: the generated episodes are stored in a buffer as (s, a, r, s') quadruples, and for training a transition is sampled uniformly;
- separate target network θ: to prevent instabilities, a separate target network is used to estimate state-value function. The weights are not trained, but copied from the trained network once in a while;
- reward clipping of $(r+\gamma \max_{a'} Q(s',a';ar{ heta}) Q(s,a; heta))$ to [-1,1].

NPFL122, Lecture 6

DDQN

DistRL

Deep Q Networks Hyperparameters



Hyperparameter	Value
minibatch size	32
replay buffer size	1M
target network update frequency	10k
discount factor	0.99
training frames	50M
RMSProp learning rate and momentum	0.00025, 0.95
initial $arepsilon$, final $arepsilon$ and frame of final $arepsilon$	1.0, 0.1, 1M
replay start size	50k
no-op max	30

Duelling

Refresh

Rainbow

Ú FAL

There have been many suggested improvements to the DQN architecture. In the end of 2017, the *Rainbow: Combining Improvements in Deep Reinforcement Learning* paper combines 7 of them into a single architecture they call *Rainbow*.



Figure 1 of the paper "Rainbow: Combining Improvements in Deep Reinforcement Learning" by Matteo Hessel et al.

NoisyNets

DistRL

Duelling

NPFL122, Lecture 6

DQN

Refresh

DDQN

PriRep



Double Q-learning

Similarly to double Q-learning, instead of

$$r+\gamma \max_{a'} Q(s',a';ar{ heta}) - Q(s,a; heta),$$

we minimize

$$r+\gamma Q(s',rgmax_{a'}Q(s',a'; heta);ar{ heta})-Q(s,a; heta).$$



Figure 1: The orange bars show the bias in a single Qlearning update when the action values are $Q(s, a) = V_*(s) + \epsilon_a$ and the errors $\{\epsilon_a\}_{a=1}^m$ are independent standard normal random variables. The second set of action values Q', used for the blue bars, was generated identically and independently. All bars are the average of 100 repetitions.

Figure 1 of the paper "Deep Reinforcement Learning with Double Q-learning" by Hado van Hasselt et al.

NPFL122, Lecture 6

DQN DDQN

Refresh

PriRep

Duelling NoisyNets

DistRL

Double Q-learning True value and an estimate All estimates and max Bias as function of state Average error 22 $\max_a Q_t(s,a) - \max_a Q_*(s,a)$ $\max_a Q_t(s,a)$ +0.61 $Q_*(s,a)$ 1 0 0 0 -0.02Double-Q estimate s.-1-2-2 $\max_{a} Q_t(s,a) - \max_{a} Q_*(s,a)$ $\max_a Q_t(s,a)$ +0.472 $Q_*(s,a)$ 2 $Q_t(s,a)$ +0.02Double-Q estimate -1 0 0 4 $\max_a Q_i$ $Q_t(s,a)$ $\max_a Q_t(s,a)$ $\max_a \ell$ 222+3.350 0 -0.02Double-Q estimate $P_*(s, a)$ $\mathbf{2}$ 26 $\mathbf{2}$ 6 6 4 state state state Figure 2 of the paper "Deep Reinforcement Learning with Double Q-learning" by Hado van Hasselt et al.

NPFL122, Lecture 6

DQN

Refresh

DDQN

PriRep D

Duelling Noisy

NoisyNets DistRL Rainbow

10/37

Double Q-learning Space Invaders Time Pilot Alien Zaxxon 2.5Value estimates -20 8 DQN estimate 2.0151.5Double DQN estimate A Phil Depoint 2 1.0Double DQN true value DQN true value 0 50 100 150 200 50 100 150 200 50 100 150 200 0 $50 \ 100 \ 150 \ 200$ 0 0 0 Training steps (in millions) Wizard of Wor Asterix Value estimates 100 80 $(log \ scale)$ DQN 401020DQN 10Double DQN Double DQN 5 0 50100 1502000 50100 150200 Wizard of Wor Asterix 4000 Double DQN 6000 Double DQN 3000 Score 4000 2000 2000 1000 DON 0 0 100 100 150150200 200 500 500 Training steps (in millions) Training steps (in millions)

Figure 3 of the paper "Deep Reinforcement Learning with Double Q-learning" by Hado van Hasselt et al.

NPFL122, Lecture 6

Refresh DQN

PriRep

DDQN

Duelling NoisyNets

DistRL

Rainbow

11/37



Double Q-learning

	DQN	Double DQN
Median	93.5%	114.7%
Mean	241.1%	330.3%

Table 1 of the paper "Deep Reinforcement Learning with Double Q-learning" by Hado van Hasselt et al.

	DQN	Double DQN	Double DQN (tuned)
Median	47.5%	88.4%	116.7%
Mean	122.0%	273.1%	475.2%

Table 2 of the paper "Deep Reinforcement Learning with Double Q-learning" by Hado van Hasselt et al.

NPFL122, Lecture 6

DDQN

Prioritized Replay

Instead of sampling the transitions uniformly from the replay buffer, we instead prefer those with a large TD error. Therefore, we sample transitions according to their probability

$$p_t \propto \left| r + \gamma \max_{a'} Q(s',a';ar{ heta}) - Q(s,a; heta)
ight|^{\omega},$$

where ω controls the shape of the distribution (which is uniform for $\omega = 0$ and corresponds to TD error for $\omega = 1$).

New transitions are inserted into the replay buffer with maximum probability to support exploration of all encountered transitions.

Refresh

PriRep Duelling

Ú F_AL

Prioritized Replay

Because we now sample transitions according to p_t instead of uniformly, on-policy distribution and sampling distribution differ. To compensate, we therefore utilize importance sampling with ratio

$$ho_t = \left(rac{1/N}{p_t}
ight)^eta$$
 .

The authors utilize in fact "for stability reasons"

$$ho_t/\max_i
ho_i.$$

NPFL122, Lecture 6

Refresh DQN

PriRep

DDQN

Prioritized Replay

Algorithm 1 Double DQN with proportional prioritization

- 1: Input: minibatch k, step-size η , replay period K and size N, exponents α and β , budget T.
- 2: Initialize replay memory $\mathcal{H} = \emptyset$, $\Delta = 0$, $p_1 = 1$
- 3: Observe S_0 and choose $A_0 \sim \pi_{\theta}(S_0)$
- 4: for t = 1 to T do
- 5: Observe S_t, R_t, γ_t
- 6: Store transition $(S_{t-1}, A_{t-1}, R_t, \gamma_t, S_t)$ in \mathcal{H} with maximal priority $p_t = \max_{i < t} p_i$
- 7: **if** $t \equiv 0 \mod K$ **then**
- 8: for j = 1 to k do

9: Sample transition
$$j \sim P(j) = p_j^{\alpha} / \sum_i p_i^{\alpha}$$

- 10: Compute importance-sampling weight $w_j = (N \cdot P(j))^{-\beta} / \max_i w_i$
- 11: Compute TD-error $\delta_j = R_j + \gamma_j Q_{\text{target}} (S_j, \arg \max_a Q(S_j, a)) Q(S_{j-1}, A_{j-1})$
- 12: Update transition priority $p_j \leftarrow |\delta_j|$
- 13: Accumulate weight-change $\Delta \leftarrow \Delta + w_j \cdot \delta_j \cdot \nabla_{\theta} Q(S_{j-1}, A_{j-1})$
- 14: **end for**
- 15: Update weights $\theta \leftarrow \theta + \eta \cdot \Delta$, reset $\Delta = 0$
- 16: From time to time copy weights into target network $\theta_{\text{target}} \leftarrow \theta$
- 17: **end if**
- 18: Choose action $A_t \sim \pi_{\theta}(S_t)$
- 19: **end for**

Refresh

Algorithm 1 of the paper "Prioritized Experience Replay" by Tom Schaul et al.

Rainbow

DistRL

NPFL122, Lecture 6

Ú F_ÁL

Duelling Networks

Instead of computing directly $Q(s, a; \theta)$, we compose it from the following quantities:

- value function for a given state s,
- advantage function computing an *advantage* of using action *a* in state *s*.

$$Q(s,a) \stackrel{\scriptscriptstyle{ ext{def}}}{=} V(f(s;\zeta);\eta) + A(f(s;\zeta),a;\psi) - rac{\sum_{a'\in\mathcal{A}}A(f(s;\zeta),a';\psi)}{|\mathcal{A}|}$$



Figure 1 of the paper "Dueling Network Architectures for Deep Reinforcement Learning" by Ziyu Wang et al.

NPFL122, Lecture 6

DQN

Refresh

DDQN

PriRep Duelling

NoisyNets

Duelling Networks



Figure 3. (a) The corridor environment. The star marks the starting state. The redness of a state signifies the reward the agent receives upon arrival. The game terminates upon reaching either reward state. The agent's actions are going up, down, left, right and no action. Plots (b), (c) and (d) shows squared error for policy evaluation with 5, 10, and 20 actions on a log-log s ale. The dueling network (Duel) consistently outperforms a conventional single-stream network (Single), with the performance gap increasing with the number of actions.

Figure 3 of the paper "Dueling Network Architectures for Deep Reinforcement Learning" by Ziyu Wang et al.

DDQN

Refresh

Ú F_AL

Duelling Networks



Figure 2 of the paper "Dueling Network Architectures for Deep Reinforcement Learning" by Ziyu Wang et al.

Duelling



Refresh

DQN

DDQN

PriRep

NoisyNets

DistRL

Rainbow

18/37



Duelling Networks

	30 no-ops		Human Starts	
	Mean	Median	Mean	Median
Prior. Duel Clip	591.9%	172.1%	567.0%	115.3%
Prior. Single	434.6%	123.7%	386.7%	112.9%
Duel Clip	373.1%	151.5%	343.8%	117.1%
Single Clip	341.2%	132.6%	302.8%	114.1%
Single	307.3%	117.8%	332.9%	110.9%
Nature DQN	227.9%	79.1%	219.6%	68.5%

Table 1 of the paper "Dueling Network Architectures for Deep Reinforcement Learning" by Ziyu Wang et al.

NPFL122, Lecture 6



Multi-step Learning

Instead of Q-learning, we use n-step variant Q-learning (to be exact, we use n-step Expected Sarsa) to maximize

$$\sum_{i=1}^n \gamma^{i-1}r_i + \gamma^n \max_{a'} Q(s',a';ar{ heta}) - Q(s,a; heta),$$

This changes the off-policy algorithm to on-policy, but it is not discussed in any way by the authors.

DDQN

Noisy Nets

Noisy Nets are neural networks whose weights and biases are perturbed by a parametric function of a noise.

The parameters $oldsymbol{ heta}$ are represented as

$$\boldsymbol{ heta} \stackrel{ ext{def}}{=} \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{arepsilon},$$

where $\boldsymbol{\varepsilon}$ is zero-mean noise with fixed statistics. We therefore learn the parameters $\boldsymbol{\zeta} \stackrel{\text{\tiny def}}{=} (\boldsymbol{\mu}, \boldsymbol{\sigma})$. Therefore, a fully connected layer

$$y = wx + b$$

is represented in the following way in Noisy Nets:

$$oldsymbol{y} = (oldsymbol{\mu}_w + oldsymbol{\sigma}_w \odot oldsymbol{arepsilon}_w)oldsymbol{x} + (oldsymbol{\mu}_b + oldsymbol{\sigma}_b \odot oldsymbol{arepsilon}_b).$$

Duelling

NPFL122, Lecture 6

DQN DDQN

Refresh

PriRep

NoisyNets DistRL



Noisy Nets

The noise ε can be for example independent Gaussian noise. However, for performance reasons, factorized Gaussian noise is used to generate a matrix of noise. If $\varepsilon_{i,j}$ is noise corresponding to a layer with i inputs and j outputs, we generate independent noise ε_i for input neurons, independent noise ε_j for output neurons, and set

$$arepsilon_{i,j} = f(arepsilon_i) f(arepsilon_j)$$

for $f(x) = \operatorname{sign}(x) \sqrt{|x|}$.

The authors generate noise samples for every batch, sharing the noise for all batch instances.

Deep Q Networks

When training a DQN, ε -greedy is no longer used and all policies are greedy, and all fully connected layers are parametrized as noisy nets.

Noisy Nets

	Bas	seline	Noi	syNet	Improvement
	Mean	Median	Mean	Median	(On median)
DQN	319	83	379	123	48%
Dueling	524	132	633	172	30%
A3Č	293	80	347	94	18%

Table 1 of the paper "Noisy Networks for Exploration" by Meire Fortunato et al.



Figure 2 of the paper "Noisy Networks for Exploration" by Meire Fortunato et al.

NPFL122, Lecture 6

Refresh DQN

DDQN

PriRep Duelling

Ú FÂL

Noisy Nets



Figure 3: Comparison of the learning curves of the average noise parameter $\overline{\Sigma}$ across five Atari games in NoisyNet-DQN. The results are averaged across 3 seeds and error bars (+/- standard deviation) are plotted.

Figure 3 of the paper "Noisy Networks for Exploration" by Meire Fortunato et al.

NPFL122, Lecture 6

Refresh DQN

PriRep Duelling

DDQN



Distributional RL

Instead of an expected return Q(s, a), we could estimate distribution of expected returns Z(s, a).

These distributions satisfy a distributional Bellman equation:

$$Z(s,a)=R(s,a)+\gamma Z(s',a').$$

The authors of the paper prove similar properties of the distributional Bellman operator compared to the regular Bellman operator, mainly being a contraction under a suitable metric (Wasserstein metric).

Duelling

DDQN

Ú F_AL

Distributional RL

The distribution of returns is modeled as a discrete distribution parametrized by the number of atoms $N \in \mathbb{N}$ and by $V_{\text{MIN}}, V_{\text{MAX}} \in \mathbb{R}$. Support of the distribution are atoms

$$\{z_i \stackrel{ ext{def}}{=} V_{ ext{MIN}} + i\Delta z: 0 \leq i < N \} \; \; ext{for} \; \Delta z \stackrel{ ext{def}}{=} rac{V_{ ext{MAX}} - V_{ ext{MIN}}}{N-1}$$

The atom probabilities are predicted using a softmax distribution as

$$Z_{oldsymbol{ heta}}(s,a) = \left\{ z_i ext{ with probability } p_i = rac{e^{f_i(s,a)}}{\sum_j e^{f_j(s,a)}}
ight\}.$$

NPFL122, Lecture 6

Refresh DQN

PriRep

DDQN

Distributional RL

After the Bellman update, the support of the distribution $R(s,a) + \gamma Z(s',a')$ is not the same as the original support. We therefore project it to the original support by proportionally mapping each atom of the Bellman update to immediate neighbors in the original support.



Figure 1 of the paper "A Distributional Perspective on Reinforcement Learning" by Marc G. Bellemare et al.

1

$$\Phiig(R(s,a)+\gamma Z(s',a')ig)_i \stackrel{ ext{def}}{=} \sum_{j=1}^N \left[1-rac{\left|[r+\gamma z_j]_{V_{ ext{MIN}}}^{V_{ ext{MAX}}}-z_i
ight|}{\Delta z}
ight]_0^r p_j(s',a').$$

The network is trained to minimize the Kullbeck-Leibler divergence between the current distribution and the (mapped) distribution of the one-step update

$$D_{ ext{KL}}ig(\Phi(R+\max_{a'}Z(s',a')||Z(s,a)ig).$$

Duelling

NPFL122, Lecture 6

Refresh

Distributional RL

Algorithm 1 Categorical Algorithm

input A transition $x_t, a_t, r_t, x_{t+1}, \gamma_t \in [0, 1]$ $Q(x_{t+1}, a) := \sum_{i} z_i p_i(x_{t+1}, a)$ $a^* \leftarrow \arg \max_a Q(x_{t+1}, a)$ $m_i = 0, \quad i \in 0, \dots, N-1$ for $j \in 0, ..., N - 1$ do # Compute the projection of $\mathcal{T}z_i$ onto the support $\{z_i\}$ $\hat{\mathcal{T}}z_j \leftarrow [r_t + \gamma_t z_j]_{V_{\text{max}}}^{V_{\text{max}}}$ $b_i \leftarrow (\mathcal{T} z_i - V_{\text{MIN}}) / \Delta z \ \# b_i \in [0, N-1]$ $l \leftarrow |b_i|, u \leftarrow [b_i]$ # Distribute probability of $\hat{\mathcal{T}}z_i$ $m_l \leftarrow m_l + p_i(x_{t+1}, a^*)(u - b_i)$ $m_u \leftarrow m_u + p_i(x_{t+1}, a^*)(b_i - l)$ end for **output** $-\sum_{i} m_i \log p_i(x_t, a_t)$ # Cross-entropy loss

Algorithm 1 of the paper "A Distributional Perspective on Reinforcement Learning" by Marc G. Bellemare et al.

Duelling

NPFL122, Lecture 6

Refresh DQN

PriRep

DDQN



Refresh

DQN

Distributional RL

	Mean	Median	> H.B.	>DQN
DQN	228%	79%	24	0
DDQN	307%	118%	33	43
DUEL.	373%	151%	37	50
Prior.	434%	124%	39	48
PR. DUEL.	592%	172%	39	44
C51	701%	178%	40	50

Figure 6 of the paper "A Distributional Perspective on Reinforcement Learning" by Marc G. Bellemare et al.



Figure 4. Learned value distribution during an episode of SPACE INVADERS. Different actions are shaded different colours. Returns below 0 (which do not occur in SPACE INVADERS) are not shown here as the agent assigns virtually no probability to them.

Figure 4 of the paper "A Distributional Perspective on Reinforcement Learning" by Marc G. Bellemare et al.

NPFL122, Lecture 6

DDQN PriRep





Distributional RL



Figure 18. SPACE INVADERS: Top-Left: Multi-modal distribution with high uncertainty. Top-Right: Subsequent frame, a more certain demise. Bottom-Left: Clear difference between actions. Bottom-Middle: Uncertain survival. Bottom-Right: Certain success. *Figure 18 of the paper "A Distributional Perspective on Reinforcement Learning" by Marc G. Bellemare et al.*

NPFL122, Lecture 6

Refresh DQN

DDQN

PriRep [



Distributional RL



Training Frames (millions)

Figure 3. Categorical DQN: Varying number of atoms in the discrete distribution. Scores are moving averages over 5 million frames. Figure 3 of the paper "A Distributional Perspective on Reinforcement Learning" by Marc G. Bellemare et al.

NPFL122, Lecture 6

Refresh DQN

PriRep

DDQN

Rainbow Architecture



Rainbow combines all described DQN extensions. Instead of 1-step updates, n-step updates are utilized, and KL divergence of the current and target return distribution is minimized:

$$D_{ ext{KL}}ig(\Phi(G_{t:t+n}+\gamma^n\max_{a'}Z(s',a'))||Z(s,a)ig).$$

The prioritized replay chooses transitions according to the probability

$$p_t \propto \Big(D_{ ext{KL}}ig(\Phi(G_{t:t+n} + \gamma^n \max_{a'} Z(s',a')) || Z(s,a) ig) \Big)^{\omega}.$$

Network utilizes duelling architecture feeding the shared representation $f(s; \zeta)$ into value computation $V(f(s; \zeta); \eta)$ and advantage computation $A_i(f(s; \zeta), a; \psi)$ for atom z_i , and the final probability of atom z_i in state s and action a is computed as

$$p_i(s,a) \stackrel{ ext{def}}{=} rac{e^{V(f(s;\zeta);\eta) + A_i(f(s;\zeta),a;\psi) - \sum_{a'\in\mathcal{A}}A_i(f(s;\zeta),a';\psi)/|\mathcal{A}|}}{\sum_j e^{V(f(s;\zeta);\eta) + A_j(f(s;\zeta),a;\psi) - \sum_{a'\in\mathcal{A}}A_j(f(s;\zeta),a';\psi)/|\mathcal{A}|}}.$$

NPFL122, Lecture 6

Refresh

DQN

Rainbow Hyperparameters

Ú	FAL

Parameter	Value
Min history to start learning	80K frames
Adam learning rate	0.0000625
Exploration ϵ	0.0
Noisy Nets σ_0	0.5
Target Network Period	32K frames
Adam ϵ	1.5×10^{-4}
Prioritization type	proportional
Prioritization exponent ω	0.5
Prioritization importance sampling β	$0.4 \rightarrow 1.0$
Multi-step returns n	3
Distributional atoms	51
Distributional min/max values	[-10, 10]

Table 1 of the paper "Rainbow: Combining Improvements in Deep Reinforcement Learning" by Matteo Hessel et al.

NPFL122, Lecture 6

Refresh DQN

DDQN PriRep

Rainbow Results





Agent	no-ops	human starts
DQN	79%	68%
DDQN (*)	117%	110%
Prioritized DDQN (*)	140%	128%
Dueling DDQN (*)	151%	117%
A3C (*)	-	116%
Noisy DQN	118%	102%
Distributional DQN	164%	125%
Rainbow	223%	153%

Table 2 of the paper "Rainbow: Combining Improvements in Deep Reinforcement Learning" byMatteo Hessel et al.



NPFL122, Lecture 6

Refresh DQN

PriRep

Duelling

DDQN

NoisyNets DistRL

stRL Rainbow

Rainbow Results





Figure 1 of the paper "Rainbow: Combining Improvements in Deep Reinforcement Learning" by Figure 3 of the paper "Rainbow: Combining Improvements in Deep Reinforcement Learning" by Matteo Hessel et al.

Duelling

NPFL122, Lecture 6

DQN DDQN

Refresh

PriRep

NoisyNets DistRL

Rainbow

35/37

Rainbow Ablations



Figure 2: Each plot shows, for several agents, the number of games where they have achieved at least a given fraction of human performance, as a function of time. From left to right we consider the 20%, 50%, 100%, 200% and 500% thresholds. On the first row we compare Rainbow to the baselines. On the second row we compare Rainbow to its ablations.

Duelling

Figure 2 of the paper "Rainbow: Combining Improvements in Deep Reinforcement Learning" by Matteo Hessel et al.

Rainbow

Refresh

DQN

DDQN PriRep

NoisyNets DistRL

Rainbow Ablations



NPFL122, Lecture 6

DQN DDQN

Refresh

N PriRep

Duelling NoisyNets

DistRL Rainbow