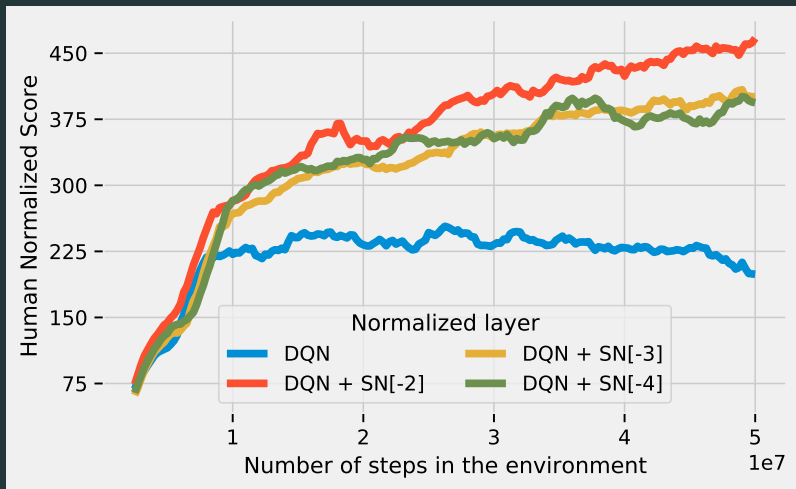


# Spectral Normalisation for Deep Reinforcement Learning

An Optimisation Perspective

---

**Florin Gogianu, Tudor Berariu, Mihaela Roşca,**  
Claudia Clopath, Lucian Buşoniu, Răzvan Paşcanu



Mean Human Normalized Score on 54 Atari games.

1. Does **smoothness** explain the performance?
2. Are there other **optimisation**-related effects at play?

## Spectral Normalization

Linear case:  $f(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ . Then:

$$f \text{ is } K\text{-Lipschitz in } \|\cdot\|_2 \Leftrightarrow \|\mathbf{W}\|_2 \leq K,$$

where  $\|\mathbf{W}\|_2$  is the largest singular value or **spectral radius** of  $\mathbf{W}$ .

## Spectral Normalization

Linear case:  $f(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$ . Then:

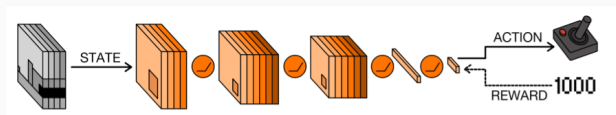
$$f \text{ is } K\text{-Lipschitz in } \|\cdot\|_2 \Leftrightarrow \|\mathbf{W}\|_2 \leq K,$$

where  $\|\mathbf{W}\|_2$  is the largest singular value or **spectral radius** of  $\mathbf{W}$ .

Approximate the largest singular value  $\rho$  of  $\mathbf{W}$  and normalize [Miyato et al., 2018]:

$$\begin{aligned}\rho &= \text{one-step-power-iteration}(\mathbf{W}_t) \\ \hat{\mathbf{W}}_t &= \mathbf{W}_t / \rho\end{aligned}$$

# Lipschitz constant of a Neural Network



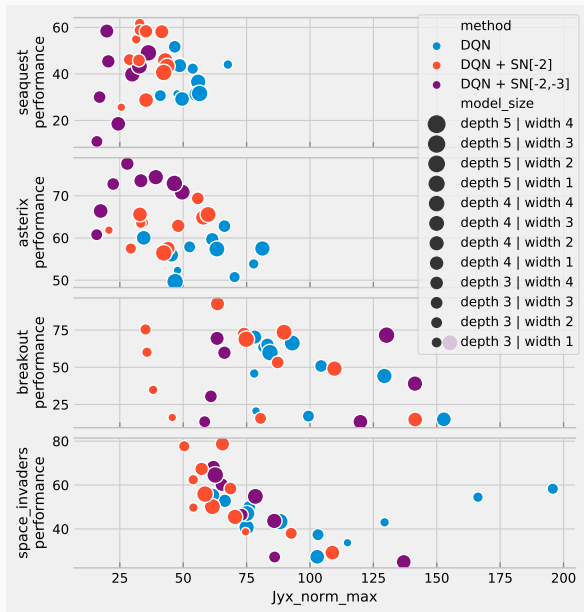
**Figure 1:** Typical architecture in DRL when learning from pixels, [Arulkumaran et al., 2017]

- Linear layers  $\phi_{fc}(\mathbf{x}) = \mathbf{W}\mathbf{x} + \mathbf{b}$  can be  $K$ -Lipschitz if we constrain the spectral radius of  $\mathbf{W}$ .
- Convolutional operators are also linear maps.
- ReLU non-linearities and Max-Pooling are 1-Lipschitz.

For a full neural network:

$$L(f) \leq \prod_{i=1} L(\phi_i)$$

# Is it the smoothness?



Smoothness of the neural network (small norm of the Jacobian  $\|J_{yx}\|$ ) is not consistently correlated with performance.

Normalised score on four MinAtar games, ten seeds each.

## The optimization perspective

MLP

$$\frac{\partial \mathcal{L}}{\partial \mathbf{W}_i} = \mathbf{J}_i \delta_L \mathbf{a}_{i-1}^\top$$

$$\frac{\partial \mathcal{L}}{\partial \mathbf{b}_i} = \mathbf{J}_i \delta_L$$

SN+bias scaling

$$\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{W}_i} = \rho^{-1} \mathbf{J}_i \hat{\delta}_L \mathbf{a}_{i-1}^\top$$

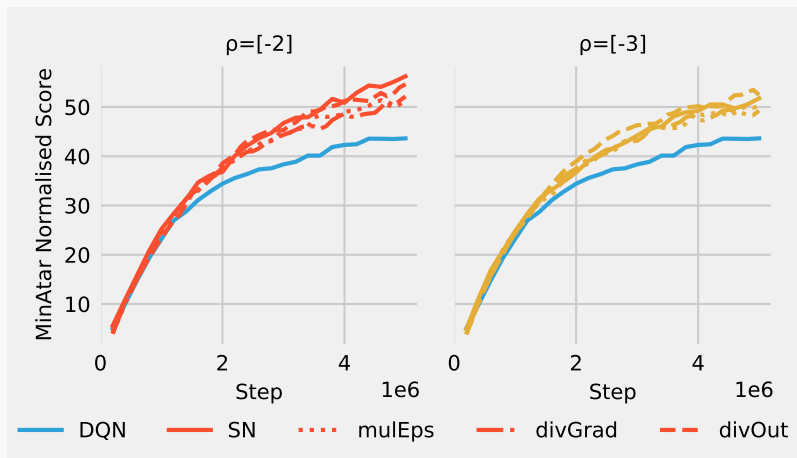
$$\frac{\partial \hat{\mathcal{L}}}{\partial \mathbf{b}_i} = \rho^{-1} \mathbf{J}_i \hat{\delta}_L$$

where  $\rho^{-1} = \prod_{i \in \mathcal{S}} \rho_i^{-1}$ . Several equivalent schedulers become apparent:

1. DivOut: output divided by  $\rho^{-1}$ .
2. DivGrad: gradient divided by  $\rho^{-1}$ .
3. MulEps: Adam's  $\epsilon$  multiplied by  $\rho^{-1}$ .



## Schedulers recover most of the effect of SN



**Figure 2:** Spectral schedulers recover SN performance. Average normalised scores over MinAtar games and four different models

# Atari results

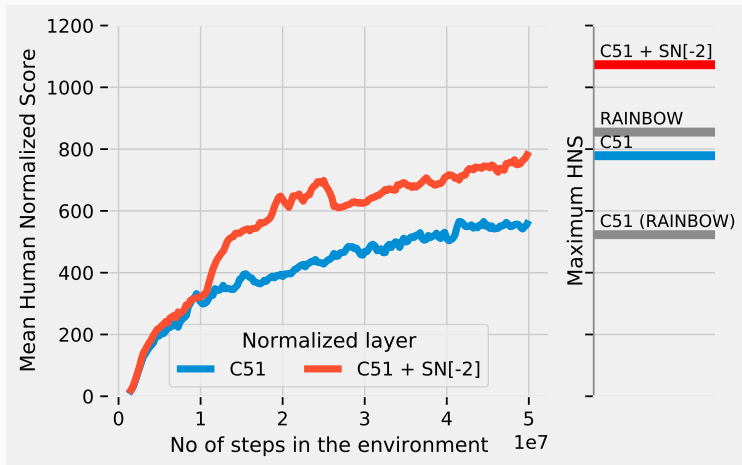




Figure 3: Categorical-DQN Performance. Mean Human Normalized Score 54 Atari games.

We improve on RAINBOW, a much more complex agent that combines many DRL advances, while using only its cost function.

**There is much to gain by designing better adapting optimisers for Deep Reinforcement Learning.**

-  Arulkumaran, K., Deisenroth, M. P., Brundage, M., and Bharath, A. A. (2017).  
**A brief survey of deep reinforcement learning.**  
*CoRR*, abs/1708.05866.
-  Miyato, T., Kataoka, T., Koyama, M., and Yoshida, Y. (2018).  
**Spectral normalization for generative adversarial networks.**