# REINFORCEMENT LEARNING
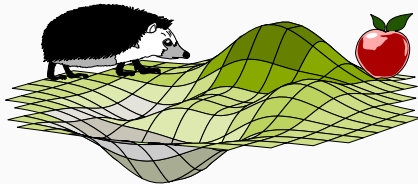
## DEEP REINFORCEMENT LEARNING

Pavel Osinenko

Consider the following $\infty$-horizon optimal control problem:

$$\max_{\kappa} J(x_0 \mid \kappa) = \int_0^\infty \rho(x(t), \kappa(x(t))) \, dt$$

$$\text{s.t.} \quad \dot{x} = f(x, u)$$

Recall the HJB:

$$\max_{u \in \mathcal{U}} \left\{ \mathcal{L}_{f(x,u)} V(x) + \rho(x, u) \right\} = 0, \quad \forall x,$$

where the Lie derivative reads:

$$\mathcal{L}_{f(x,u)} V(x) = \nabla V(x)^{\mathsf{T}} f(x, u).$$

We consider here a deterministic environment only for the ease of exposition

The corresponding Hamiltonian reads:

$$\mathcal{H}(x, u \mid h) := \mathcal{L}_{f(x,u)} h(x) + p(x,u)$$

for a generic function $h$ of $x$.

If $\kappa^*$ is the optimal policy, then the HJB can be rewritten as:

$$\mathcal{H}(x, \kappa^*(x) \mid V) = 0 , \quad \forall x$$

We will learn the value function via a deep neural network assuming a suitable (stabilizing and exploring) behavior policy generating $u(t), t \geq 0$
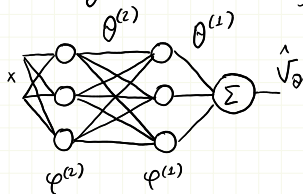
For the sake of brevity, let's take a two-layer network, whereas the argument will work for any topology in a similar manner.
So,

$$\forall x \quad \hat{V}_\theta(x) := \theta^{(1)^T} \varphi^{(1)} \left( \theta^{(2)^T} \varphi^{(2)}(x) \right),$$

where $\theta^{(1)}$ is the weight vector of the output layer, $\theta^{(2)}$ is the weight matrix of the hidden layer, $\varphi^{(1)}, \varphi^{(2)}$ are the corresponding activation functions.

Details: $\theta^{(1)}$ is $N_{c_1} \times 1$
$\theta^{(2)}$ is $N_{c_1} \times N_{c_2}$

Let's write down the gradient:

$$\nabla_x \hat{V}_\theta(x) = \left( \nabla_x \varphi^{(2)}(x) \, \theta_j^{(2)} \right)_j^{N_{c1}} \left( \frac{\partial \varphi^{(1)}}{\partial O^{(2)}} \right)^T \theta^{(1)} \, ,$$

where $O^{(2)}$ is the output of the hidden layer:

$$O^{(2)} := \theta^{(2)^T} \varphi^{(2)}(x) \, .$$

The notation $\theta_j^{(2)}$ means the $j$th row of the matrix $\theta^{(2)}$ and $\left( \bullet_j \right)_j^{N_{c1}}$ means horizontal stacking over $j$, i.e.,

$$\left( \bullet_1, \bullet_2, \ldots, \bullet_{N_{c1}} \right)$$

The Hamiltonian for the critic reads:

$$\mathcal{H}(x, u \mid \hat{V}_\theta) = \mathcal{L}_{f(x,u)} \hat{V}_\theta(x) + \rho(x, u)$$

$$= \nabla_x \hat{V}_\theta(x)^T f(x, u) + \rho(x, u)$$

By a universal approximation theorem,

$$V(x) = \theta^{*(1)^T} \varphi^{(1)}\left( \theta^{*(2)^T} \varphi^{(2)}(x) \right) + \delta(x)$$

(on a compact).

Here, $\theta^* = \{ \theta^{(1)^T}, \theta^{(2)^T} \}$ is the set of the best weights

Accordingly, the Hamiltonian approximation error reads:

$$S_{\mathcal{H}}(x, u) := \mathcal{H}(x, u \mid V) - \mathcal{H}(x, u \mid \hat{V}_{\theta^*})$$

$$= \mathcal{L}_{f(x,u)} V(x) - \mathcal{L}_{f(x,u)} \hat{V}_{\theta^*}$$

$$= \mathcal{L}_{f(x,u)} \delta(x) .$$

The Hamiltonian TD in turn is:

$$e_{\mathcal{H}}(\theta \mid x, u) := \mathcal{H}(x, u \mid \hat{V}_{\theta}) - \mathcal{H}(x, \kappa^*(x) \mid V)$$

$$= \mathcal{H}(x, u \mid \hat{V}_{\theta}) ,$$

where $\kappa^*$ is the optimal policy

In more detail,

$$H(x, u | \hat{V}_\theta) = \theta^{(1)T} \frac{\partial \varphi^{(1)}}{\partial O^{(1)}} \left( \left( \nabla_x \varphi^{(2)}(x) \, \theta_j^{(2)} \right)_j^{N_{C1}} \right)^T f(x, u) + \rho(x, u).$$

Let us compute the gradients of the Hamiltonian TD:

$$\nabla_{\theta^{(1)}} H(x, u | \hat{V}_\theta) = \frac{\partial \varphi^{(1)}}{\partial O^{(1)}} \left( \left( \nabla_x \varphi^{(2)}(x) \, \theta_j^{(2)} \right)_j^{N_{C1}} \right)^T f(x, u)$$

(this part resembles the one in the case of a 1-layer network)

For the second part, let's introduce the notation

$$G(\theta^{(2)}, x) := \frac{\partial \varphi^{(2)}}{\partial O^{(2)}} \left( \left( \nabla_x \varphi^{(2)}(x) \theta_j^{(2)} \right)_j^{N_{c_L}} \right)^T.$$

Then, $\mathcal{H}(x, u \mid \hat{V}_\theta) = \theta^{(2)^T} G(\theta^{(2)}, x) f(x, u) + \rho(x, u)$.

Let's do some unwrapping:

$$\mathcal{H}(x, u \mid \hat{V}_\theta) = \left( \left( \theta^{(2)^T} G_\ell^T(\theta^{(2)}, x) \right)_\ell^n \right)^T f(x, u) + \rho(x, u)$$

$$= \sum_{\ell=1}^n \theta^{(2)^T} G_\ell^T(\theta^{(2)}, x) f_\ell(x, u) + \rho(x, u)$$

Let's abuse the notation and interpret $\theta^{(2)}$ as a vector while computing

$$\nabla_{\theta^{(2)}}$$

to mean, actually, $\nabla_{vec(\theta^{(2)})}$ where $vec(\cdot)$ is a transformation of a matrix into a vector. Alternatively, we could always interpret $\theta^{(2)}$ as a vector and reshape it into matrix to compute $\hat{V}_\theta$.

Then,

$$\nabla_{\theta^{(2)}} \, \mathcal{H}(x, u \,|\, \hat{V}_\theta) = \sum_{l=1}^{n} f_l(x, u) \, \nabla_{\theta^{(2)}} \, G_l^{\top}(\theta^{(2)}, x) \, \theta^{(2)^{\top}}$$

So far, we have :

$$\nabla_\theta \mathcal{H}\left(x, u \,\middle|\, \hat{V}_\theta\right) = \begin{pmatrix} \frac{\partial \varphi^{(1)}}{\partial \theta^{(2)}} \left( \left( \nabla_x \varphi^{(2)}(x)\, \theta_j^{(2)} \right)_j^{N_{c1}} \right)^T f(x, u) \\[4ex] \sum_{l=1}^{n} f_l(x, u)\, \nabla_{\theta^{(2)}} \left( \left( \left( \nabla_x \varphi^{(2)}(x)\, \theta_j^{(2)} \right)_j^{N_{c1}} \frac{\partial \varphi^{(1)}}{\partial \theta^{(2)}} \right) \theta^{(1)^T} \right)_l \end{pmatrix} .$$

Be aware that $\frac{\partial \varphi^{(1)}}{\partial \theta^{(2)}}$ in general depends nonlinearly on $\theta^{(2)}$ !

Let us linearize $\mathcal{H}(x, u \mid \hat{V}_{\theta^*})$ around $\theta$ :

$$\mathcal{H}(x, u \mid \hat{V}_{\theta^*}) = \theta^{(1)^T} G(\theta^{(2)}, x) f(x, u) \; -$$

$$\left( \begin{array}{c} \dfrac{\partial \varphi^{(1)}}{\partial \theta^{(2)}} \left( \left( \nabla_x \varphi^{(2)}(x) \, \theta_j^{(2)} \right)_j^{N_{c1}} \right)^T f(x, u) \\[20pt] \sum_{\ell=1}^{n} f_\ell(x, u) \, \nabla_{\theta^{(2)}} \left( \left( \nabla_x \varphi^{(2)}(x) \, \theta_j^{(2)} \right)_j^{N_{c1}} \dfrac{\partial \varphi^{(1)}}{\partial \theta^{(2)}} \right) \theta^{(1)^T} \right)_\ell \end{array} \right)^T \Bigg|_{\theta = \theta} \operatorname{vec}(\tilde{\theta}) \; +$$

$$O\left( \| \operatorname{vec}(\tilde{\theta}) \|^2 \right) + \rho(x, u) ,$$

where $\operatorname{vec}(\tilde{\theta}) := \operatorname{vec}(\theta) - \operatorname{vec}(\theta^*) = \left( \theta^{(1)^T}, \operatorname{vec}(\theta^{(2)^T})^T \right)^T - \left( \theta^{*(1)^T}, \operatorname{vec}(\theta^{*(2)^T})^T \right)^T$

(cont.)

Expressing $\rho(x,u)$ from the last equality and plugging into $e_{\mathcal{H}}(\theta|x,u)$ yields:

$$e_{\mathcal{H}}(\theta|x,u) = \left( \begin{array}{c} \dfrac{\partial \varphi^{(i)}}{\partial \mathcal{O}^{(i)}} \left( \left( \nabla_x \varphi^{(i)}(x)\, \theta_j^{(i)} \right)_j^{N_{ci}} \right)^T f(x,u) \\[2em] \displaystyle\sum_{l=1}^{n} f_l(x,u)\, \nabla_{\theta^{(i)}} \left( \left( \nabla_x \varphi^{(i)}(x)\, \theta_j^{(i)} \right)_j^{N_{ci}} \dfrac{\partial \varphi^{(i)}}{\partial \mathcal{O}^{(i)}} \right) \theta^{(i)T} \end{array} \right)_l \Bigg|^T_{\theta=\hat{\theta}} \; vec(\tilde{\theta}) \; +$$

$$\mathcal{H}(x,u \,|\, \hat{V}_{\theta^*}) \; + \; \mathcal{O}\left( \| vec(\tilde{\theta}) \|^2 \right)$$

Let the data vector be defined as:

$$w(x,u,\theta) := \left\| \begin{array}{c} \dfrac{\partial \varphi^{(2)}}{\partial \theta^{(2)}} \left( \left( \nabla_x \varphi^{(2)}(x)\, \theta_j^{(2)} \right)_j^{N_{c2}} \right)^T f(x,u) \\[20pt] \sum_{l=1}^{n} f_l(x,u)\, \nabla_{\theta^{(2)}} \left( \left( \nabla_x \varphi^{(2)}(x)\, \theta_j^{(2)} \right)_j^{N_{c2}} \dfrac{\partial \varphi^{(2)}}{\partial \theta^{(2)}} \right)_l \theta^{(2)T} \end{array} \right\|_{\theta = \theta} .$$

Notice it depends on $\theta$, but we do not expect troubles with persistence of excitation because of that since the weights shouldn't normally converge to zero

Summarizing, we have:

$$e_{\mathcal{H}}(\theta/x,u) = vec(\tilde{\theta})^T \varpi(x,u,\theta) + \mathcal{H}(x,u|\hat{V}_{\theta^*}) +$$

$$O(\|vec(\tilde{\theta})\|^2).$$

Relating this to the true Hamiltonian gives:

$$e_{\mathcal{H}}(\theta/x,u) = vec(\tilde{\theta})^T \varpi(x,u,\theta) + \mathcal{H}(x,u|V) -$$

$$\mathcal{S}_{\mathcal{H}}(x,u) + O(\|vec(\tilde{\theta})\|^2)$$

As in the case of a shallow critic, let's define the loss as :

$$J^c\left(\theta \mid \{x(t_k), u(t_k)\}_k^M\right) := \frac{1}{2} \sum_{k=1}^{M} \frac{e_{\mathcal{H}}^2\left(\theta \mid x(t_k), u(t_k)\right)}{\left(w_k^T w_k + 1\right)^2},$$

where $M$ is the size of the experience replay $\{x(t_k), u(t_k)\}_k^M$ with $t_M = t$. The data points may be equally distributed. Also, we used a short-hand notation

$$w_k := w\left(x(t_k), u(t_k), \theta\right)$$

Accordingly, let's update the critic weights by gradient descent over mini-batches:

$$\frac{d}{dt} vec(\Theta) := -\alpha \nabla_\Theta J^c\left(\Theta \,\middle|\, \{x(t_k), u(t_k)\}_k^M\right):$$

$$= -\alpha \sum_{k=1}^M \frac{e_{\mathcal{H}}\left(\Theta \,\middle|\, x(t_k), u(t_k)\right) \nabla_\Theta e_{\mathcal{H}}\left(\Theta \,\middle|\, x(t_k), u(t_k)\right)}{\left(\omega_k^T \omega_k + 1\right)^2} ,$$

where $\alpha$ is the learning rate.

Since $\nabla_\Theta e_{\mathcal{H}}\left(\Theta \,\middle|\, x(t_k), u(t_k)\right) = \omega_k$ , we have:

$$\frac{d}{dt} vec(\Theta) = -\alpha \sum_{k=1}^M \frac{e_{\mathcal{H}}\left(\Theta \,\middle|\, x(t_k), u(t_k)\right) \omega_k}{\left(\omega_k^T \omega_k + 1\right)^2}$$

Next, recalling that

$$e_{\mathcal{H}}(\theta | x, u) = \text{vec}(\tilde{\theta})^T \omega(x, u, \theta) + \mathcal{H}(x, u | V) - \delta_{\mathcal{H}}(x, u) + \mathcal{O}(\|\text{vec}(\tilde{\theta})\|^2),$$

rewrite the update rule as follows:

$$\frac{d}{dt} \text{vec}(\theta) = -\alpha \sum_{k=1}^{M} \frac{\omega_k \, \omega_k^T}{\left(\omega_k^T \omega_k + 1\right)^2} \text{vec}(\tilde{\theta}) +$$

$$\alpha \sum_{k=1}^{M} \frac{\omega_k \left(\mathcal{H}(x(t_k), u(t_k) | V) - \delta_{\mathcal{H}}(x(t_k), u(t_k)) + \mathcal{O}(\|\text{vec}(\tilde{\theta})\|^2)\right)}{\left(\omega_k^T \omega_k + 1\right)^2}$$

Define the following data matrix :

$$\mathcal{E}(t) := \sum_{k=1}^{M} \frac{\omega_k \, \omega_k^{\top}}{\left(\omega_k^{\top} \omega_k + 1\right)^2} \; .$$

Then, proceeding to the norm and noticing that $\frac{d}{dt} vec(\theta) = \frac{d}{dt} vec(\tilde{\theta})$, we have :

$$\frac{d}{dt} \|vec(\tilde{\theta})\|^2 = -\alpha \, \|vec(\tilde{\theta})\|^2_{\mathcal{E}(t)} \; +$$

$$\alpha \sum_{k=1}^{M} \frac{\omega_k \left( vec(\tilde{\theta})^{\top} \mathcal{H} \left( x(t_k), u(t_k) \,|\, V \right) - \mathcal{S}_{\mathcal{H}}\left( x(t_k), u(t_k) \right) + \mathcal{O}\left( \|vec(\tilde{\theta})\|^3 \right) \right)}{\left( \omega_k^{\top} \omega_k + 1 \right)^2} \; ,$$

where $\|vec(\tilde{\theta})\|^2_{\mathcal{E}(t)} = vec(\tilde{\theta})^{\top} \mathcal{E}(t) \, vec(\tilde{\theta})$

Since the behaviour policy is stabilizing, noticing that $\frac{\|\omega\|}{(\|\omega\|^2 + 1)^2} \leq \frac{1}{2}, \forall \omega$, and assuming $\mathcal{H}, V$ are continuous, there are constants $\bar{\mathcal{H}}, \bar{\delta}_{\mathcal{H}}, C > 0$ s.t.

$$\alpha \sum_{k=1}^{M} \frac{\omega_k \left( \text{vec}(\tilde{\theta})^\top \mathcal{H}(x(t_k), u(t_k) | V) - \delta_{\mathcal{H}}(x(t_k), u(t_k)) + C \|\text{vec}(\tilde{\theta})\|^3 \right)}{\left( \omega_k^\top \omega_k + 1 \right)^2} \leq$$

$$\frac{\alpha M}{2} \left( \|\text{vec}(\tilde{\theta})\| \cdot \left( M + \bar{\delta}_{\mathcal{H}} \right) + C \|\text{vec}(\tilde{\theta})\|^3 \right)$$

Assuming persistence of excitation, i.e., $\exists \varepsilon > 0$ s.t. $\forall t \geq 0 \;\; \mathcal{E}(t) \succcurlyeq \varepsilon I$, we may claim

$$\frac{d}{dt} \| \text{vec}(\tilde{\theta}) \|^2 \leq -\alpha \varepsilon \| \text{vec}(\tilde{\theta}) \|^2 +$$

$$\frac{\alpha M}{2} \left( \| \text{vec}(\tilde{\theta}) \| \cdot \left( M + \bar{\delta}_{\mathcal{H}} \right) + C \, \| \text{vec}(\tilde{\theta}) \|^3 \right).$$

In absense of the linearization error $C \| \text{vec}(\tilde{\theta}) \|^3$, the argumentation about the critic weight convergence would have been the same as in the case of a shallow network: the weights converge from everywhere into a vicinity of the best weights

$$\frac{d}{dt} \left\| vec(\tilde{\theta}) \right\|^2 \leqslant -\alpha \varepsilon \left\| vec(\tilde{\theta}) \right\|^2 +$$

$$\frac{\alpha M}{2} \left( \left\| vec(\tilde{\theta}) \right\| \cdot \left( M + \bar{\delta}_{\mathcal{H}} \right) + C \left\| vec(\tilde{\theta}) \right\|^3 \right)$$

This vicinity depends primarily on two factors:

1. how good the chosen network is capable of approximating $V$ and thus $\mathcal{H}$ (affects $\bar{\delta}_{\mathcal{H}}$);

2. how close the behavior policy is to the optimal one (affects $M$)

$$\frac{d}{dt} \| vec(\tilde{\theta}) \|^2 \leq -\alpha \varepsilon \| vec(\tilde{\theta}) \|^2 +$$

$$\frac{\alpha M}{2} \left( \| vec(\tilde{\theta}) \| \cdot \left( M + \bar{\delta}_{\mathcal{H}} \right) + C \| vec(\tilde{\theta}) \|^3 \right)$$

(cont.)  The  point  2.  could  in  theory  be
tackled  by  combination  of  exploration/exploitation
phases  thus,  hopefully,  approaching  the
optimal  policy.

Question:  is  such  a  method  on-  or  off-policy?

Important remark:  there  is  no  guarantee  of
stabilizability  during  the  exploitation  following  this
recipe

$$\frac{d}{dt} \| vec(\tilde{\theta}) \|^2 \leq - \alpha \varepsilon \| vec(\tilde{\theta}) \|^2 +$$

$$\frac{\alpha M}{2} \left( \| vec(\tilde{\theta}) \| \cdot \left( M + \bar{\delta}_H \right) + C \| vec(\tilde{\theta}) \|^3 \right)$$

(cont.) The key principle utilized in the convergence analysis is that the first, quadratic term dominates the linear term outside a vicinity of zero. But in the case of a deep neural network critic, there is a higher-order term which makes an essential difference

$$\frac{d}{dt} \| vec(\tilde{\theta}) \|^2 \leqslant -\alpha \varepsilon \| vec(\tilde{\theta}) \|^2 +$$

$$\frac{\alpha M}{2} \left( \| vec(\tilde{\theta}) \| \cdot \left( M + \bar{\delta}_{\mathcal{H}} \right) + C \| vec(\tilde{\theta}) \|^3 \right)$$

(cont.)   The principle here is exactly the opposite: if the weight error is not too large initially, it converges. Now, we are faced by the problem that the convergence domain is limited by this factor and also by what we discussed earlier — the effect of the Hamiltonian approximation error and goodness of the behavior policy

Convergence of the actor weights is performed in a similar manner. Again, linearization errors affect the convergence.

But local character of convergence is a general problem of deep learning.

A further remark on the derivation of the update rule : in general, back propagation is the systematic routine to analytically derive the gradients

$$\nabla_\theta \mathcal{H}(x, u \mid \hat{V}_\theta)$$

(cont.)

Furthermore, modifications of the gradient descent usual in deep learning may be applied, such as momentum techniques, Nesterov accelerated gradient descent; gradient clipping to avoid fading/explosion (at deep layers) may also be of use, as well as such techniques as dropout (switching off the network partly). Next, ReLu activation functions can be employed whose gradients are either one or zero

An overview of deep RL methods:

- DQN ("deep QL") and double DQN

- DDPG ("DPG + DQN")

- TD3 (twin delayed deep deterministic): a modification of DDPG

- αGo Zero: uses a single deep network to fit the policy and value, updates via Monte-Carlo tree search