

# Bound on the Weights of a Neural Network Under Nesterov's Accelerated Gradient Flow

The following is a derivation of a bound on the weights of a vanilla neural network under Nesterov's Accelerated Gradient Flow (the continuous version of Nesterov's Accelerated Gradient Descent). Finding such a bound was stated as an "open question" in Gluch and Urbanke's "Noether" paper (<https://arxiv.org/abs/2104.05508>). The following derivation assumes familiarity with their paper:

We start with "conservation law" for Nesterov's Accelerated Gradient Flow:

$$\left\langle W^{(h)}, \ddot{W}^{(h)} + \frac{3}{t}\dot{W}^{(h)} \right\rangle - \left\langle W^{(h+1)}, \ddot{W}^{(h+1)} + \frac{3}{t}\dot{W}^{(h+1)} \right\rangle = 0 \quad (1)$$

First multiply by  $t$

$$t \left( \left\langle W^{(h)}, \ddot{W}^{(h)} \right\rangle - \left\langle W^{(h+1)}, \ddot{W}^{(h+1)} \right\rangle \right) + 3 \left( \left\langle W^{(h)}, \dot{W}^{(h)} \right\rangle - \left\langle W^{(h+1)}, \dot{W}^{(h+1)} \right\rangle \right) = 0 \quad (2)$$

Notice that the second term can be expressed as a derivative of  $\|W^{(h)}\|_F^2 - \|W^{(h+1)}\|_F^2$ :

$$t \left( \left\langle W^{(h)}, \ddot{W}^{(h)} \right\rangle - \left\langle W^{(h+1)}, \ddot{W}^{(h+1)} \right\rangle \right) + \frac{3}{2} \frac{d}{dt} \left( \|W^{(h)}\|_F^2 - \|W^{(h+1)}\|_F^2 \right) = 0 \quad (3)$$

Next, we can perform integration by parts on the first term (showing just layer (h) for simplicity).

$$\int \left\langle W^{(h)}, \ddot{W}^{(h)} \right\rangle t \, dt \quad (4)$$

$$= \int \sum_{i,j} \left( W_{ij}^{(h)} t \right) \ddot{W}_{ij}^{(h)} \, dt \quad (5)$$

$$= \sum_{i,j} \int \left( W_{ij}^{(h)} t \right) \ddot{W}_{ij}^{(h)} \, dt \quad (6)$$

$$= \sum_{i,j} \left[ - \int \left( \dot{W}_{ij}^{(h)} t + W_{ij}^{(h)} \right) \dot{W}_{ij}^{(h)} \, dt + W_{ij}^{(h)} \dot{W}_{ij}^{(h)} t \right] \quad (7)$$

$$= \left\langle W^{(h)}, \dot{W}^{(h)} \right\rangle t - \frac{\|W^{(h)}\|^2}{2} - \int \|\dot{W}^{(h)}\|^2 t \, dt \quad (8)$$

Hence, equation (3) becomes

$$\frac{d}{dt} \left( \|W^{(h)}\|_F^2 - \|W^{(h+1)}\|_F^2 \right) + \frac{d}{dt} \left( \left\langle W^{(h)}, \dot{W}^{(h)} \right\rangle t - \left\langle W^{(h+1)}, \dot{W}^{(h+1)} \right\rangle t \right) = t \left( \|\dot{W}^{(h)}\|^2 - \|\dot{W}^{(h+1)}\|^2 \right) \quad (9)$$

Now, note that

$$\left\langle W^{(h)}, \dot{W}^{(h)} \right\rangle t = \frac{d}{dt} \left( \frac{1}{2} \|W^{(h)}\|^2 t \right) - \frac{1}{2} \|W^{(h)}\|^2 \quad (10)$$

Therefore,

$$\frac{d^2}{dt^2} \left( t \|W^{(h)}\|_F^2 - \|W^{(h+1)}\|_F^2 \right) + \frac{d}{dt} \left( \|W^{(h)}\|_F^2 - \|W^{(h+1)}\|_F^2 \right) = 2 t \left( \|\dot{W}^{(h)}\|^2 - \|\dot{W}^{(h+1)}\|^2 \right) \quad (11)$$

For simplicity, let  $\alpha = \|W^{(h)}\|_F^2 - \|W^{(h+1)}\|_F^2$ . Then,

$$\frac{d^2}{dt^2} (t\alpha) + \dot{\alpha} = 2 t \left( \|\dot{W}^{(h)}\|^2 - \|\dot{W}^{(h+1)}\|^2 \right) \quad (12)$$

$$\ddot{\alpha} + \frac{3}{t}\dot{\alpha} = 2 \left( \|\dot{W}^{(h)}\|^2 - \|\dot{W}^{(h+1)}\|^2 \right) \quad (13)$$

Next, acknowledge that

$$\ddot{\alpha} + \frac{3}{t}\dot{\alpha} \leq 2 \left( \|\dot{W}^{(h)}\|^2 + \|\dot{W}^{(h+1)}\|^2 \right) \leq 4 \left( \frac{1}{2} \left( \|\dot{W}^{(h)}\|^2 + \|\dot{W}^{(h+1)}\|^2 \right) + L(\omega) - L(\omega^*) \right) \quad (14)$$

When summed across all layers, the quantity on the right is simply the Hamiltonian of the system. Moreover, the Hamiltonian is decreasing,

$$\dot{\mathcal{H}} = -\frac{3}{t} \|\dot{\omega}\|^2 \quad (15)$$

It follows, then, that  $\mathcal{H} \leq \mathcal{H}_0$ . Hence,

$$\ddot{\alpha} + \frac{3}{t}\dot{\alpha} \leq 4\mathcal{H}_0 \quad (16)$$

$$\ddot{\alpha} \leq -\frac{3}{t}\dot{\alpha} + 4\mathcal{H}_0 \quad (17)$$

Here, we can make the change of variables:

$$\beta = \dot{\alpha} - \mathcal{H}_o t \quad (18)$$

$$\dot{\beta} = \ddot{\alpha} - \mathcal{H}_o \quad (19)$$

This gives

$$\dot{\beta} \leq -\frac{3}{t}\beta \quad (20)$$

We can now apply Gronwall's Inequality, which yields

$$\beta \leq \beta(t_o) e^{\int_{t_o}^t (-\frac{3}{t}) dt} \quad (21)$$

$$\beta \leq \beta(t_o) \left(\frac{t_o}{t}\right)^3 \quad (22)$$

If we set  $t_o = 0$ , then

$$\dot{\alpha} - \mathcal{H}_o t \leq 0 \quad (23)$$

$$\alpha(t) - \alpha(0) \leq \frac{1}{2} \mathcal{H}_o t^2 \quad (24)$$

$$\left[ \sum_{h=1}^{K-1} \left| \|W^{(h)}(t)\|_F^2 - \|W^{(h+1)}(t)\|_F^2 \right| \right] - \left[ \sum_{h=1}^{K-1} \left| \|W^{(h)}(0)\|_F^2 - \|W^{(h+1)}(0)\|_F^2 \right| \right] \leq \frac{1}{2} \mathcal{H}_o t^2 \quad (25)$$

Or, if we note that  $\mathcal{H}_o = L(\omega(0)) - L^*$  when the velocities are initialized at 0, we arrive at

$$\left[ \sum_{h=1}^{K-1} \left| \|W^{(h)}(t)\|_F^2 - \|W^{(h+1)}(t)\|_F^2 \right| \right] - \left[ \sum_{h=1}^{K-1} \left| \|W^{(h)}(0)\|_F^2 - \|W^{(h+1)}(0)\|_F^2 \right| \right] \leq \frac{1}{2} (L(\omega(0)) - L^*) t^2 \quad (26)$$

which looks a lot like the bound that was derived for the case of Newtonian dynamics. The absolute value comes in by observing the symmetry of the problem (e.g. by defining  $\alpha$  by its negative and thus achieving a lower bound on the original  $\alpha$ ).