

In the last six weeks we In September I developed a framework for auto-differentiation of arbitrary neural networks. This framework enables us to minimize losses of the form  $\min g(\nabla f(x))$ . Using this framework we ran a series of experiments aimed at capturing a physically based solution to  $\min_{\phi} \sum_{t=1}^T \|\dot{\mathbf{r}}_t \cdot \nabla \phi(\mathbf{r}_t)\|_2^2$ . Early results produced trivial solutions e.g.  $\nabla \phi(\mathbf{r}) = 0 \implies \phi(\mathbf{r}_t) = c$ . A gradient normalization term was added to reduce the likelihood of these results, and for numeric stability the cosine distance between  $\dot{\mathbf{r}}_t$  and  $\nabla \phi(\mathbf{r}_t)$  replaced the dot product of the two terms:

$$\min_{\phi} \sum_{t=1}^T \left( \left| \frac{\dot{\mathbf{r}}_t \cdot \nabla \phi(\mathbf{r}_t)}{\|\dot{\mathbf{r}}_t\|_2^2 * \|\nabla \phi(\mathbf{r}_t)\|_2^2} \right| + \left(1 - \|\nabla \phi(\mathbf{r}_t)\|_2^2\right)^2 \right)$$

This objective was trained by drawing 4000 samples uniformly from trajectories of each planet resulting in a mini-batch of  $n = 32000$  samples. Multiple fully connected base networks were trained with 1-3 layers and 4-256 hidden units per layer. Results were stable, with a small single layer network easily minimizing the objective. When evaluating  $\phi(\mathbf{r})$  however, the value would drift by a scalar factor resulting in saturation or vanishing  $\phi(\mathbf{r})$  thus the scale of  $\phi(\mathbf{r})$  was fixed by adding a second normalization term,  $(0.5 - \frac{1}{n} \sum \phi(\mathbf{r}))^2$ .

This gave a strong distinction between the value of  $\phi$  for each of the planets at the beginning of training, however the constant value of  $\phi(\mathbf{r})$  along each planet's trajectories converged to the fixed scale, 0.5, during training. In addition the gradient normalization term,  $(1 - \|\nabla \phi(\mathbf{r}_t)\|_2^2)^2$  accounted for a large portion of the final loss leading to a relaxation of this constraint to a modified hinge loss with zero cost for  $1 \leq \|\nabla \phi(\mathbf{r}_t)\| \leq 2$

$$\min_{\phi} \sum_{n=1}^N \left( \left| \frac{\dot{\mathbf{r}}_t \cdot \nabla \phi(\mathbf{r}_n)}{\|\dot{\mathbf{r}}_t\|_2^2 * \|\nabla \phi(\mathbf{r}_n)\|_2^2} \right| + \max\left(\|\nabla \phi(\mathbf{r}_n)\|_2^2 - 2, 0\right) + \max\left(1 - \|\nabla \phi(\mathbf{r}_n)\|_2^2, 0\right) \right)$$

This objective, while still minimizing the dot-product,  $\dot{\mathbf{r}}_t \cdot \nabla \phi(\mathbf{r}_n)$  resulted in stable  $\phi(\mathbf{r})$  unique to each planet's trajectory that remained separable during training. With large networks (three fully connected layers with 50+ units) these planetary constants would still converge, however for single layer networks (50, 12, 4 hidden units tested)  $\phi()$  would remain reasonably stable conditioned on each planet and the variance of a mini-batch was not decreasing as it did with larger networks.

Despite this the constants learned show no obvious relation to the energy function of the orbits they encode. Some issues that I will address in the coming weeks are the discontinuity between sampled points in trajectories, exploring a smoothness function over  $\phi(\mathbf{r})$  and evaluating  $\phi(\mathbf{r})$  over the entire domain throughout training.