

On penalty terms

Giulio Galvan

28 settembre 2015

We have seen that:

$$\frac{\partial g}{\partial W^{rec}} = \sum_{t=1}^T \frac{\partial g_t}{\partial \mathbf{a}^t} \cdot \sum_{k=1}^t \frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} \cdot \frac{\partial \mathbf{a}^k}{\partial W^{rec}}, \quad (1)$$

where

$$\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} = \prod_{i=t-1}^k \text{diag}(\sigma'(\mathbf{a}^i)) \cdot W^{rec}. \quad (2)$$

tends to vanish.

The idea is to develop a penalty term which express a preference for solutions where the components $\frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k}$ are far from zero, hence, learning a model which exhibit long memory.

A first attempt could be:

$$\Gamma \triangleq \sum_{t=1}^T \sum_{k=1}^t \frac{1}{\left\| \frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} \right\|^2}. \quad (3)$$

Γ treats all temporal steps equally: we could assign more importance to distant temporal steps, which are the most critical, modifying *Gamma* as follows:

$$\Gamma \triangleq \sum_{t=1}^T \sum_{k=1}^t \frac{\sigma(t-k)}{\left\| \frac{\partial \mathbf{a}^t}{\partial \mathbf{a}^k} \right\|^2}, \quad (4)$$

where $\sigma(\cdot)$ assign different weights depending on the temporal distance $t - k$, for example $\sigma(h) = \exp\{h\}$