

To understand the effect of the spectral radius, consider the simple case of back-propagation with a Jacobian matrix \mathbf{J} that does not change with t . This case happens, for example, when the network is purely linear. Suppose that \mathbf{J} has an eigenvector \mathbf{v} with corresponding eigenvalue λ . Consider what happens as we propagate a gradient vector backwards through time. If we begin with a gradient vector \mathbf{g} , then after one step of back-propagation, we will have $\mathbf{J}\mathbf{g}$, and after n steps we will have $\mathbf{J}^n\mathbf{g}$. Now consider what happens if we instead back-propagate a perturbed version of \mathbf{g} . If we begin with $\mathbf{g} + \delta\mathbf{v}$, then after one step, we will have $\mathbf{J}(\mathbf{g} + \delta\mathbf{v})$. After n steps, we will have $\mathbf{J}^n(\mathbf{g} + \delta\mathbf{v})$. From this we can see that back-propagation starting from \mathbf{g} and back-propagation starting from $\mathbf{g} + \delta\mathbf{v}$ diverge by $\delta\mathbf{J}^n\mathbf{v}$ after n steps of back-propagation. If \mathbf{v} is chosen to be a unit eigenvector of \mathbf{J} with eigenvalue λ , then multiplication by the Jacobian simply scales the difference at each step. The two executions of back-propagation are separated by a distance of $\delta|\lambda|^n$. When \mathbf{v} corresponds to the largest value of $|\lambda|$, this perturbation achieves the widest possible separation of an initial perturbation of size δ .

When $|\lambda| > 1$, the deviation size $\delta|\lambda|^n$ grows exponentially large. When $|\lambda| < 1$, the deviation size becomes exponentially small.

Of course, this example assumed that the Jacobian was the same at every time step, corresponding to a recurrent network with no nonlinearity. When a nonlinearity is present, the derivative of the nonlinearity will approach zero on many time steps, and help to prevent the explosion resulting from a large spectral radius. Indeed, the most recent work on echo state networks advocates using a spectral radius much larger than unity (Yildiz *et al.*, 2012; Jaeger, 2012).

Everything we have said about back-propagation via repeated matrix multiplication applies equally to forward propagation in a network with no nonlinearity, where the state $\mathbf{h}^{(t+1)} = \mathbf{h}^{(t)\top}\mathbf{W}$.

When a linear map \mathbf{W}^\top always shrinks \mathbf{h} as measured by the L^2 norm, then we say that the map is **contractive**. When the spectral radius is less than one, the mapping from $\mathbf{h}^{(t)}$ to $\mathbf{h}^{(t+1)}$ is contractive, so a small change becomes smaller after each time step. This necessarily makes the network forget information about the past when we use a finite level of precision (such as 32 bit integers) to store the state vector.

The Jacobian matrix tells us how a small change of $\mathbf{h}^{(t)}$ propagates one step forward, or equivalently, how the gradient on $\mathbf{h}^{(t+1)}$ propagates one step backward, during back-propagation. Note that neither \mathbf{W} nor \mathbf{J} need to be symmetric (although they are square and real), so they can have complex-valued eigenvalues and eigenvectors, with imaginary components corresponding to potentially oscillatory