

limited  $L^1$  norm. Usually we do not know the size of the constraint region that we impose by using weight decay with coefficient  $\alpha^*$  because the value of  $\alpha^*$  does not directly tell us the value of  $k$ . In principle, one can solve for  $k$ , but the relationship between  $k$  and  $\alpha^*$  depends on the form of  $J$ . While we do not know the exact size of the constraint region, we can control it roughly by increasing or decreasing  $\alpha$  in order to grow or shrink the constraint region. Larger  $\alpha$  will result in a smaller constraint region. Smaller  $\alpha$  will result in a larger constraint region.

Sometimes we may wish to use explicit constraints rather than penalties. As described in section 4.4, we can modify algorithms such as stochastic gradient descent to take a step downhill on  $J(\boldsymbol{\theta})$  and then project  $\boldsymbol{\theta}$  back to the nearest point that satisfies  $\Omega(\boldsymbol{\theta}) < k$ . This can be useful if we have an idea of what value of  $k$  is appropriate and do not want to spend time searching for the value of  $\alpha$  that corresponds to this  $k$ .

Another reason to use explicit constraints and reprojection rather than enforcing constraints with penalties is that penalties can cause non-convex optimization procedures to get stuck in local minima corresponding to small  $\boldsymbol{\theta}$ . When training neural networks, this usually manifests as neural networks that train with several “dead units.” These are units that do not contribute much to the behavior of the function learned by the network because the weights going into or out of them are all very small. When training with a penalty on the norm of the weights, these configurations can be locally optimal, even if it is possible to significantly reduce  $J$  by making the weights larger. Explicit constraints implemented by re-projection can work much better in these cases because they do not encourage the weights to approach the origin. Explicit constraints implemented by re-projection only have an effect when the weights become large and attempt to leave the constraint region.

Finally, explicit constraints with reprojection can be useful because they impose some stability on the optimization procedure. When using high learning rates, it is possible to encounter a positive feedback loop in which large weights induce large gradients which then induce a large update to the weights. If these updates consistently increase the size of the weights, then  $\boldsymbol{\theta}$  rapidly moves away from the origin until numerical overflow occurs. Explicit constraints with reprojection prevent this feedback loop from continuing to increase the magnitude of the weights without bound. Hinton *et al.* (2012c) recommend using constraints combined with a high learning rate to allow rapid exploration of parameter space while maintaining some stability.

In particular, Hinton *et al.* (2012c) recommend a strategy introduced by Srebro and Shraibman (2005): constraining the norm of each *column* of the weight matrix