

more than amplifying sensor noise or compression artifacts in such cases. This motivates introducing a small, positive regularization parameter  $\lambda$  to bias the estimate of the standard deviation. Alternately, one can constrain the denominator to be at least  $\epsilon$ . Given an input image  $\mathbf{X}$ , GCN produces an output image  $\mathbf{X}'$ , defined such that

$$X'_{i,j,k} = s \frac{X_{i,j,k} - \bar{X}}{\max \left\{ \epsilon, \sqrt{\lambda + \frac{1}{3rc} \sum_{i=1}^r \sum_{j=1}^c \sum_{k=1}^3 (X_{i,j,k} - \bar{X})^2} \right\}}. \quad (12.3)$$

Datasets consisting of large images cropped to interesting objects are unlikely to contain any images with nearly constant intensity. In these cases, it is safe to practically ignore the small denominator problem by setting  $\lambda = 0$  and avoid division by 0 in extremely rare cases by setting  $\epsilon$  to an extremely low value like  $10^{-8}$ . This is the approach used by Goodfellow *et al.* (2013a) on the CIFAR-10 dataset. Small images cropped randomly are more likely to have nearly constant intensity, making aggressive regularization more useful. Coates *et al.* (2011) used  $\epsilon = 0$  and  $\lambda = 10$  on small, randomly selected patches drawn from CIFAR-10.

The scale parameter  $s$  can usually be set to 1, as done by Coates *et al.* (2011), or chosen to make each individual pixel have standard deviation across examples close to 1, as done by Goodfellow *et al.* (2013a).

The standard deviation in equation 12.3 is just a rescaling of the  $L^2$  norm of the image (assuming the mean of the image has already been removed). It is preferable to define GCN in terms of standard deviation rather than  $L^2$  norm because the standard deviation includes division by the number of pixels, so GCN based on standard deviation allows the same  $s$  to be used regardless of image size. However, the observation that the  $L^2$  norm is proportional to the standard deviation can help build a useful intuition. One can understand GCN as mapping examples to a spherical shell. See figure 12.1 for an illustration. This can be a useful property because neural networks are often better at responding to directions in space rather than exact locations. Responding to multiple distances in the same direction requires hidden units with collinear weight vectors but different biases. Such coordination can be difficult for the learning algorithm to discover. Additionally, many shallow graphical models have problems with representing multiple separated modes along the same line. GCN avoids these problems by reducing each example to a direction rather than a direction and a distance.

Counterintuitively, there is a preprocessing operation known as **sphering** and it is not the same operation as GCN. Sphering does not refer to making the data lie on a spherical shell, but rather to rescaling the principal components to have