



Figure 14.4: A denoising autoencoder is trained to map a corrupted data point  $\tilde{\mathbf{x}}$  back to the original data point  $\mathbf{x}$ . We illustrate training examples  $\mathbf{x}$  as red crosses lying near a low-dimensional manifold illustrated with the bold black line. We illustrate the corruption process  $C(\tilde{\mathbf{x}} | \mathbf{x})$  with a gray circle of equiprobable corruptions. A gray arrow demonstrates how one training example is transformed into one sample from this corruption process. When the denoising autoencoder is trained to minimize the average of squared errors  $\|g(f(\tilde{\mathbf{x}})) - \mathbf{x}\|^2$ , the reconstruction  $g(f(\tilde{\mathbf{x}}))$  estimates  $\mathbb{E}_{\mathbf{x}, \tilde{\mathbf{x}} \sim p_{\text{data}}(\mathbf{x}) C(\tilde{\mathbf{x}} | \mathbf{x})}[\mathbf{x} | \tilde{\mathbf{x}}]$ . The vector  $g(f(\tilde{\mathbf{x}})) - \tilde{\mathbf{x}}$  points approximately towards the nearest point on the manifold, since  $g(f(\tilde{\mathbf{x}}))$  estimates the center of mass of the clean points  $\mathbf{x}$  which could have given rise to  $\tilde{\mathbf{x}}$ . The autoencoder thus learns a vector field  $g(f(\mathbf{x})) - \mathbf{x}$  indicated by the green arrows. This vector field estimates the score  $\nabla_{\mathbf{x}} \log p_{\text{data}}(\mathbf{x})$  up to a multiplicative factor that is the average root mean square reconstruction error.