



Figure 3.6: The KL divergence is asymmetric. Suppose we have a distribution  $p(x)$  and wish to approximate it with another distribution  $q(x)$ . We have the choice of minimizing either  $D_{\text{KL}}(p||q)$  or  $D_{\text{KL}}(q||p)$ . We illustrate the effect of this choice using a mixture of two Gaussians for  $p$ , and a single Gaussian for  $q$ . The choice of which direction of the KL divergence to use is problem dependent. Some applications require an approximation that usually places high probability anywhere that the true distribution places high probability, while other applications require an approximation that rarely places high probability anywhere that the true distribution places low probability. The choice of the direction of the KL divergence reflects which of these considerations takes priority for each application. *(Left)* The effect of minimizing  $D_{\text{KL}}(p||q)$ . In this case, we select a  $q$  that has high probability where  $p$  has high probability. When  $p$  has multiple modes,  $q$  chooses to blur the modes together, in order to put high probability mass on all of them. *(Right)* The effect of minimizing  $D_{\text{KL}}(q||p)$ . In this case, we select a  $q$  that has low probability where  $p$  has low probability. When  $p$  has multiple modes that are sufficiently widely separated, as in this figure, the KL divergence is minimized by choosing a single mode, to avoid putting probability mass in the low-probability areas between modes of  $p$ . Here, we illustrate the outcome when  $q$  is chosen to emphasize the left mode. We could also have achieved an equal value of the KL divergence by choosing the right mode. If the modes are not separated by a sufficiently strong low-probability region, then this direction of the KL divergence can still choose to blur the modes.