



Figure 20.10: An illustration of the neural autoregressive density estimator (NADE). The hidden units are organized in groups  $\mathbf{h}^{(j)}$  so that only the inputs  $x_1, \dots, x_i$  participate in computing  $\mathbf{h}^{(i)}$  and predicting  $P(x_j | x_{j-1}, \dots, x_1)$ , for  $j > i$ . NADE is differentiated from earlier neural auto-regressive networks by the use of a particular weight sharing pattern:  $W'_{j,k,i} = W_{k,i}$  is shared (indicated in the figure by the use of the same line pattern for every instance of a replicated weight) for all the weights going out from  $x_i$  to the  $k$ -th unit of any group  $j \geq i$ . Recall that the vector  $(W_{1,i}, W_{2,i}, \dots, W_{n,i})$  is denoted  $\mathbf{W}_{:,i}$ .

Larochelle and Murray (2011) chose this sharing scheme so that forward propagation in a NADE model loosely resembles the computations performed in mean field inference to fill in missing inputs in an RBM. This mean field inference corresponds to running a recurrent network with shared weights and the first step of that inference is the same as in NADE. The only difference is that with NADE, the output weights connecting the hidden units to the output are parametrized independently from the weights connecting the input units to the hidden units. In the RBM, the hidden-to-output weights are the transpose of the input-to-hidden weights. The NADE architecture can be extended to mimic not just one time step of the mean field recurrent inference but to mimic  $k$  steps. This approach is called NADE- $k$  (Raiko *et al.*, 2014).

As mentioned previously, auto-regressive networks may be extended to process continuous-valued data. A particularly powerful and generic way of parametrizing a continuous density is as a Gaussian mixture (introduced in section 3.9.6) with mixture weights  $\alpha_i$  (the coefficient or prior probability for component  $i$ ), per-component conditional mean  $\mu_i$  and per-component conditional variance  $\sigma_i^2$ . A model called RNADE (Uria *et al.*, 2013) uses this parametrization to extend NADE to real values. As with other mixture density networks, the parameters of this