

1. The parametrization of each $P(x_i | x_{i-1}, \dots, x_1)$ by a neural network with $(i - 1) \times k$ inputs and k outputs (if the variables are discrete and take k values, encoded one-hot) allows one to estimate the conditional probability without requiring an exponential number of parameters (and examples), yet still is able to capture high-order dependencies between the random variables.
2. Instead of having a different neural network for the prediction of each x_i , a *left-to-right* connectivity illustrated in figure 20.9 allows one to merge all the neural networks into one. Equivalently, it means that the hidden layer features computed for predicting x_i can be reused for predicting x_{i+k} ($k > 0$). The hidden units are thus organized in *groups* that have the particularity that all the units in the i -th group only depend on the input values x_1, \dots, x_i . The parameters used to compute these hidden units are jointly optimized to improve the prediction of all the variables in the sequence. This is an instance of the *reuse principle* that recurs throughout deep learning in scenarios ranging from recurrent and convolutional network architectures to multi-task and transfer learning.

Each $P(x_i | x_{i-1}, \dots, x_1)$ can represent a conditional distribution by having outputs of the neural network predict *parameters* of the conditional distribution of x_i , as discussed in section 6.2.1.1. Although the original neural auto-regressive networks were initially evaluated in the context of purely discrete multivariate data (with a sigmoid output for a Bernoulli variable or softmax output for a multinoulli variable) it is natural to extend such models to continuous variables or joint distributions involving both discrete and continuous variables.

20.10.10 NADE

The **neural autoregressive density estimator** (NADE) is a very successful recent form of neural auto-regressive network (Larochelle and Murray, 2011). The connectivity is the same as for the original neural auto-regressive network of Bengio and Bengio (2000b) but NADE introduces an additional parameter sharing scheme, as illustrated in figure 20.10. The parameters of the hidden units of different groups j are shared.

The weights $W'_{j,k,i}$ from the i -th input x_i to the k -th element of the j -th group of hidden unit $h_k^{(j)}$ ($j \geq i$) are shared among the groups:

$$W'_{j,k,i} = W_{k,i}. \quad (20.83)$$

The remaining weights, where $j < i$, are zero.