

set error is then chosen as having found the best hyperparameters. See the left of figure 11.2 for an illustration of a grid of hyperparameter values.

How should the lists of values to search over be chosen? In the case of numerical (ordered) hyperparameters, the smallest and largest element of each list is chosen conservatively, based on prior experience with similar experiments, to make sure that the optimal value is very likely to be in the selected range. Typically, a grid search involves picking values approximately on a *logarithmic scale*, e.g., a learning rate taken within the set  $\{.1, .01, 10^{-3}, 10^{-4}, 10^{-5}\}$ , or a number of hidden units taken with the set  $\{50, 100, 200, 500, 1000, 2000\}$ .

Grid search usually performs best when it is performed repeatedly. For example, suppose that we ran a grid search over a hyperparameter  $\alpha$  using values of  $\{-1, 0, 1\}$ . If the best value found is 1, then we underestimated the range in which the best  $\alpha$  lies and we should shift the grid and run another search with  $\alpha$  in, for example,  $\{1, 2, 3\}$ . If we find that the best value of  $\alpha$  is 0, then we may wish to refine our estimate by zooming in and running a grid search over  $\{-.1, 0, .1\}$ .

The obvious problem with grid search is that its computational cost grows exponentially with the number of hyperparameters. If there are  $m$  hyperparameters, each taking at most  $n$  values, then the number of training and evaluation trials required grows as  $O(n^m)$ . The trials may be run in parallel and exploit loose parallelism (with almost no need for communication between different machines carrying out the search) Unfortunately, due to the exponential cost of grid search, even parallelization may not provide a satisfactory size of search.

#### 11.4.4 Random Search

Fortunately, there is an alternative to grid search that is as simple to program, more convenient to use, and converges much faster to good values of the hyperparameters: random search (Bergstra and Bengio, 2012).

A random search proceeds as follows. First we define a marginal distribution for each hyperparameter, e.g., a Bernoulli or multinoulli for binary or discrete hyperparameters, or a uniform distribution on a log-scale for positive real-valued hyperparameters. For example,

$$\text{log\_learning\_rate} \sim u(-1, -5) \tag{11.2}$$

$$\text{learning\_rate} = 10^{\text{log\_learning\_rate}}. \tag{11.3}$$

where  $u(a, b)$  indicates a sample of the uniform distribution in the interval  $(a, b)$ . Similarly the `log_number_of_hidden_units` may be sampled from  $u(\log(50), \log(2000))$ .