

region containing all the points  $\mathbf{x}$  that have the same set of  $k$  nearest neighbors in the training set. For  $k = 1$ , the number of distinguishable regions cannot be more than the number of training examples.

While the  $k$ -nearest neighbors algorithm copies the output from nearby training examples, most kernel machines interpolate between training set outputs associated with nearby training examples. An important class of kernels is the family of **local kernels** where  $k(\mathbf{u}, \mathbf{v})$  is large when  $\mathbf{u} = \mathbf{v}$  and decreases as  $\mathbf{u}$  and  $\mathbf{v}$  grow farther apart from each other. A local kernel can be thought of as a similarity function that performs template matching, by measuring how closely a test example  $\mathbf{x}$  resembles each training example  $\mathbf{x}^{(i)}$ . Much of the modern motivation for deep learning is derived from studying the limitations of local template matching and how deep models are able to succeed in cases where local template matching fails (Bengio *et al.*, 2006b).

Decision trees also suffer from the limitations of exclusively smoothness-based learning because they break the input space into as many regions as there are leaves and use a separate parameter (or sometimes many parameters for extensions of decision trees) in each region. If the target function requires a tree with at least  $n$  leaves to be represented accurately, then at least  $n$  training examples are required to fit the tree. A multiple of  $n$  is needed to achieve some level of statistical confidence in the predicted output.

In general, to distinguish  $O(k)$  regions in input space, all of these methods require  $O(k)$  examples. Typically there are  $O(k)$  parameters, with  $O(1)$  parameters associated with each of the  $O(k)$  regions. The case of a nearest neighbor scenario, where each training example can be used to define at most one region, is illustrated in figure 5.10.

Is there a way to represent a complex function that has many more regions to be distinguished than the number of training examples? Clearly, assuming only smoothness of the underlying function will not allow a learner to do that. For example, imagine that the target function is a kind of checkerboard. A checkerboard contains many variations but there is a simple structure to them. Imagine what happens when the number of training examples is substantially smaller than the number of black and white squares on the checkerboard. Based on only local generalization and the smoothness or local constancy prior, we would be guaranteed to correctly guess the color of a new point if it lies within the same checkerboard square as a training example. There is no guarantee that the learner could correctly extend the checkerboard pattern to points lying in squares that do not contain training examples. With this prior alone, the only information that an example tells us is the color of its square, and the only way to get the colors of the