

as **support vectors**.

Kernel machines also suffer from a high computational cost of training when the dataset is large. We will revisit this idea in section 5.9. Kernel machines with generic kernels struggle to generalize well. We will explain why in section 5.11. The modern incarnation of deep learning was designed to overcome these limitations of kernel machines. The current deep learning renaissance began when Hinton *et al.* (2006) demonstrated that a neural network could outperform the RBF kernel SVM on the MNIST benchmark.

5.7.3 Other Simple Supervised Learning Algorithms

We have already briefly encountered another non-probabilistic supervised learning algorithm, nearest neighbor regression. More generally, k -nearest neighbors is a family of techniques that can be used for classification or regression. As a non-parametric learning algorithm, k -nearest neighbors is not restricted to a fixed number of parameters. We usually think of the k -nearest neighbors algorithm as not having any parameters, but rather implementing a simple function of the training data. In fact, there is not even really a training stage or learning process. Instead, at test time, when we want to produce an output y for a new test input \mathbf{x} , we find the k -nearest neighbors to \mathbf{x} in the training data \mathbf{X} . We then return the average of the corresponding y values in the training set. This works for essentially any kind of supervised learning where we can define an average over y values. In the case of classification, we can average over one-hot code vectors \mathbf{c} with $c_y = 1$ and $c_i = 0$ for all other values of i . We can then interpret the average over these one-hot codes as giving a probability distribution over classes. As a non-parametric learning algorithm, k -nearest neighbor can achieve very high capacity. For example, suppose we have a multiclass classification task and measure performance with 0-1 loss. In this setting, 1-nearest neighbor converges to double the Bayes error as the number of training examples approaches infinity. The error in excess of the Bayes error results from choosing a single neighbor by breaking ties between equally distant neighbors randomly. When there is infinite training data, all test points \mathbf{x} will have infinitely many training set neighbors at distance zero. If we allow the algorithm to use all of these neighbors to vote, rather than randomly choosing one of them, the procedure converges to the Bayes error rate. The high capacity of k -nearest neighbors allows it to obtain high accuracy given a large training set. However, it does so at high computational cost, and it may generalize very badly given a small, finite training set. One weakness of k -nearest neighbors is that it cannot learn that one feature is more discriminative than another. For example, imagine we have a regression task with $\mathbf{x} \in \mathbb{R}^{100}$ drawn from an isotropic Gaussian