

# **Overview of Supervised Learning**

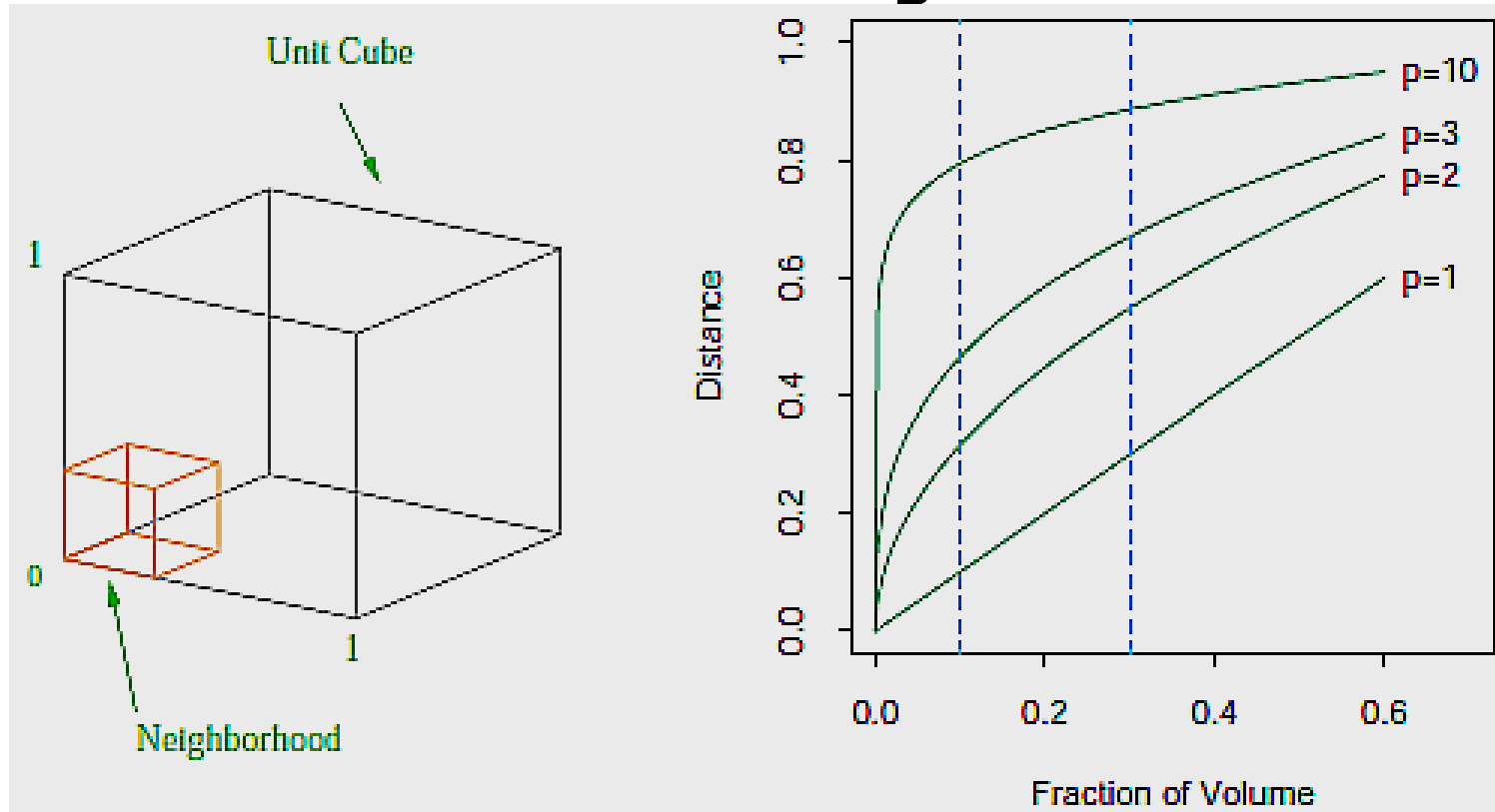
Chap – 2 ; - Part - II

**T. Hastie, R.Tibshirani, J. Friedman, "The Elements of Statistical Learning: Data Mining, Inference and Prediction", Springer Series in Statistics, 2009**

# Local Methods in High Dimensions

- Two learning techniques for prediction so far: *the stable but biased linear model* and the *less stable but apparently less biased class of  $k$ -nearest-neighbor estimates*.
- It would seem that with a reasonably large set of training data, we could always approximate the theoretically optimal conditional expectation by  $k$ -nearest-neighbor averaging, since we should be able to find a fairly large neighborhood of observations close to any  $x$  and average them.
- This approach and our intuition breaks down in high dimensions, and the phenomenon is commonly referred to as the **curse of dimensionality**.

# Local Methods in High Dimensions



**FIGURE 2.6.** The curse of dimensionality is well illustrated by a sub-cubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the sub cube needed to capture a fraction  $r$  of the volume of the data, for different dimensions  $p$ . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.

- Construct another uniform example. Suppose we have 1000 training examples  $x_i$  generated uniformly on  $[-1, 1]^p$ . Assume that the true relationship between  $X$  and  $Y$  is

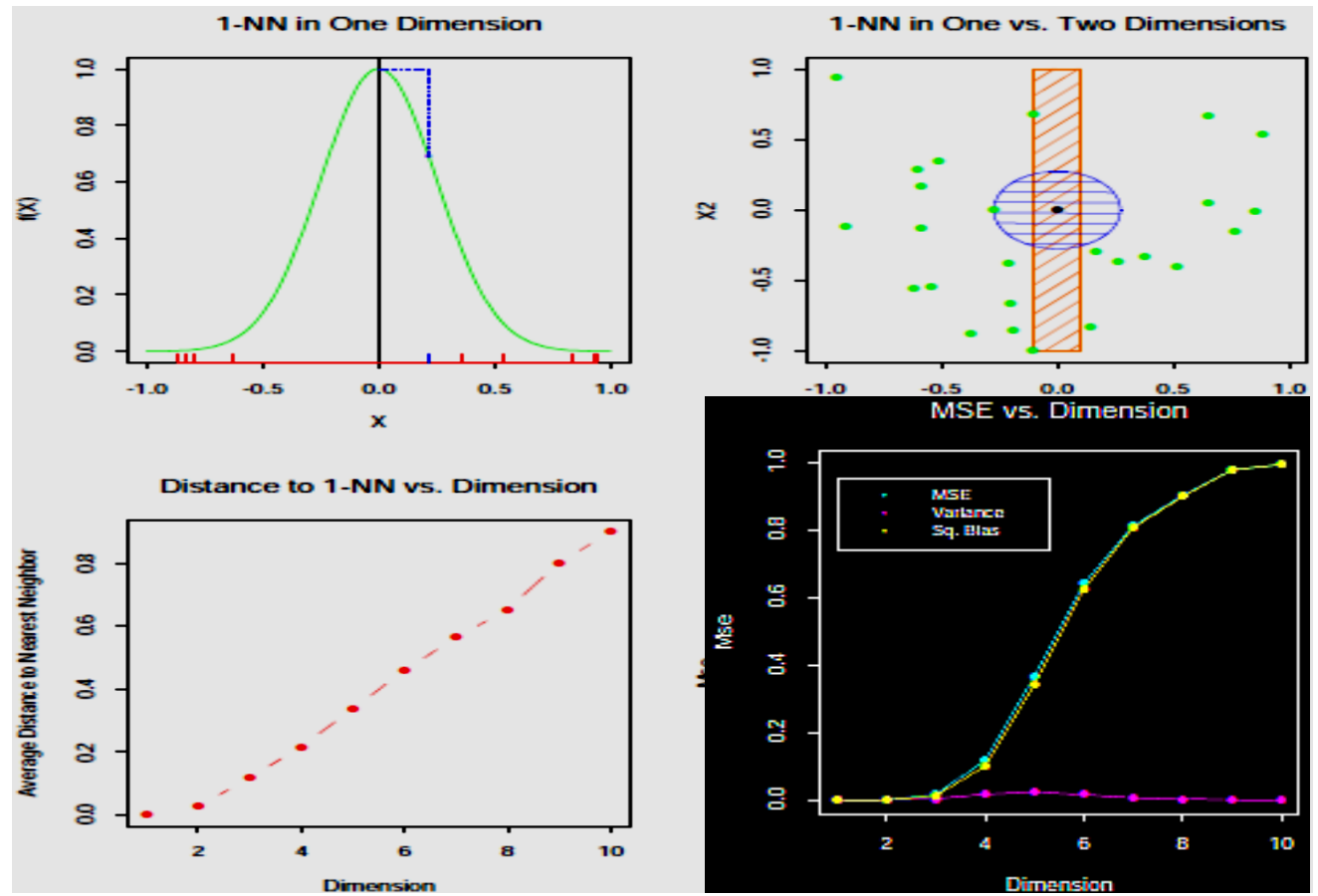
$$Y = f(X) = e^{-8\|X\|^2},$$

without any measurement error. We use the 1 –nearest-neighbor rule to predict  $y_0$  at the test-point  $x_0 = 0$ . Denote the training set by  $\tau$ .

Compute the expected prediction error at  $x_0$  for our procedure, averaging over all such samples of size 1000. The mean squared error (*MSE*) for estimating  $f(0)$ :

$$\begin{aligned} MSE(x_0) &= E_{\tau} [f(x_0) - \hat{y}_0]^2 \\ &= E_{\tau} [\hat{y}_0 - E_{\tau}(\hat{y}_0)]^2 + [E_{\tau}(\hat{y}_0) - f(x_0)]^2 \\ &= Var_{\tau}(\hat{y}_0) + Bias^2(\hat{y}_0) \end{aligned} \quad 2.25$$

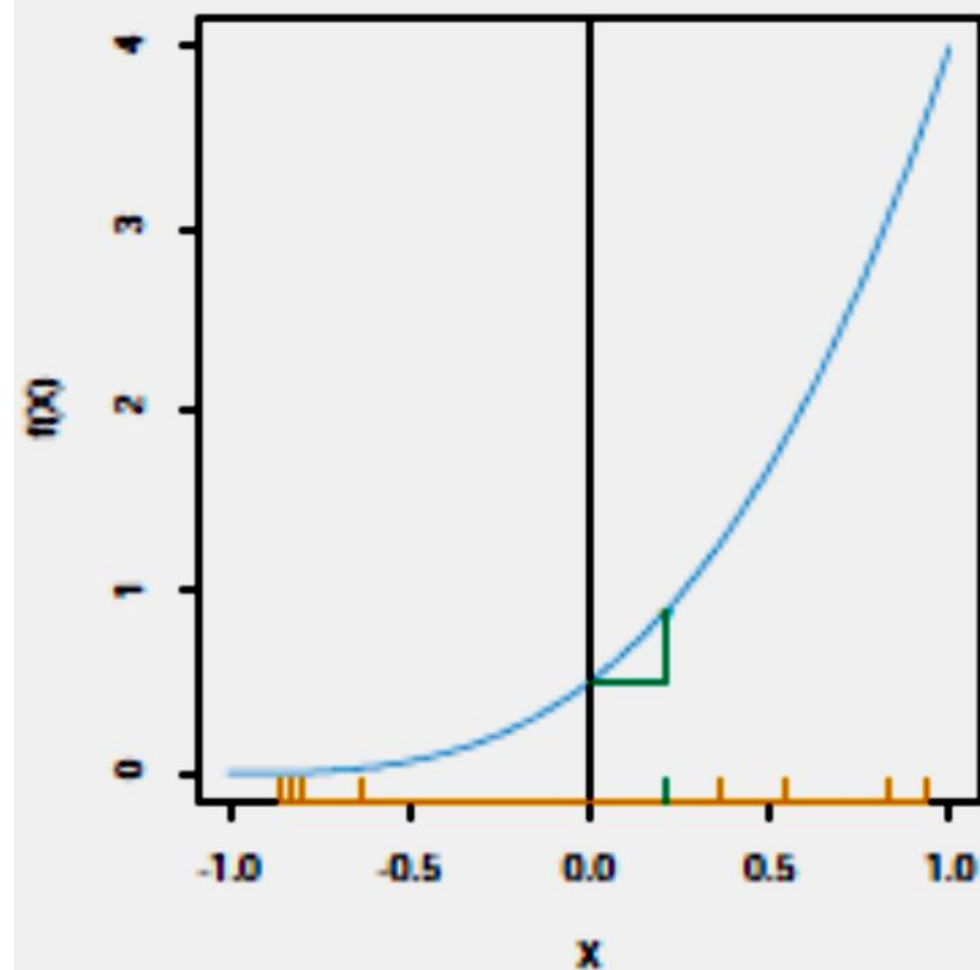
## Local Methods in High Dimensions



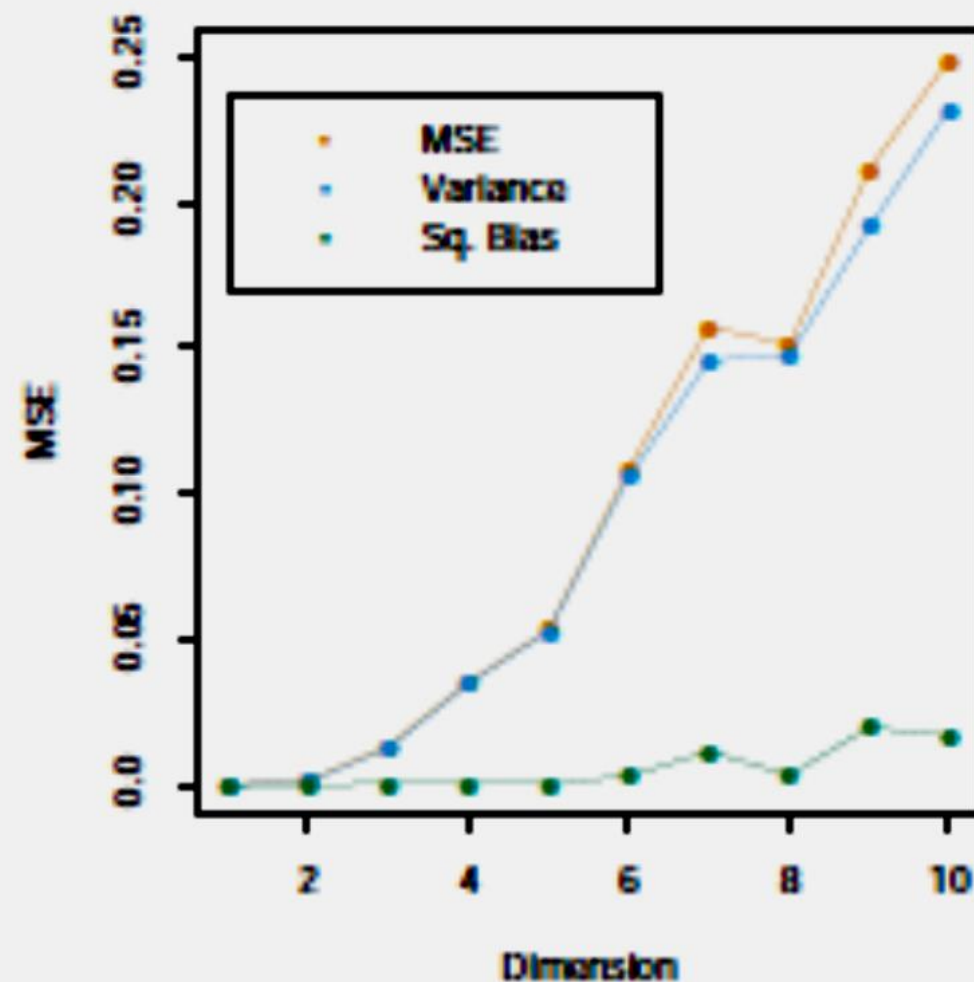
- FIGURE 2.7.** A simulation example, demonstrating the curse of dimensionality and its effect on MSE, bias and variance. The input features are uniformly distributed in  $[-1, 1]^p$  for  $p = 1, \dots, 10$ . The top left panel shows the target function (no noise) in  $R$ :  $f(X) = e^{-8\|X\|^2}$ , and demonstrates the error that 1-nearest neighbor makes in estimating  $f(0)$ . The training point is indicated by the blue tick mark. The top right panel illustrates why the radius of the 1-nearest neighborhood increases with dimension  $p$ . The lower left panel shows the average radius of the 1-nearest neighborhoods. The lower-right panel shows the MSE, squared bias and variance curves as a function of dimension  $p$ .

- **In Figure 2.7** the *MSE* is broken down into two components - variance and squared bias. Such a decomposition is always possible and often useful, and is known as the *bias–variance decomposition*.
- As the dimension increases, the nearest neighbor tends to stray further from the target point, and both bias and variance are incurred. By  $p = 10$ , for more than 99% of the samples the nearest neighbor is a distance greater than 0.5 from the origin. As  $p$  increases, the estimate tends to be 0 more often than not, and hence the *MSE* levels off at 1.0, as does the bias, and the variance starts dropping.
- Similar phenomena occur more generally. The complexity of functions of many variables can grow exponentially with the dimension, and if we wish to be able to estimate such functions with the same accuracy as function in low dimensions, then we need the size of our training set to grow exponentially as well.

1-NN in One Dimension



MSE vs. Dimension



**FIGURE 2.8.** A simulation example with the same setup as in Figure 2.7. Here the function is constant in all but one dimension:  $F(X) = \frac{1}{2}(X_1 + 1)^3$ . The variance dominates.

- If the function always involves only a few dimensions as in **Figure 2.8**, then the variance can dominate instead.
- Suppose, on the other hand, that we know that the relationship between  $Y$  and  $X$  is linear,

$$Y = X^T \beta + \varepsilon,$$

where  $\varepsilon \sim N(0, \sigma^2)$  and we fit the model by least squares to the training data. For an arbitrary test point  $x_0$ , we have  $\hat{y} = x_0^T \hat{\beta}$ , which can be written as

$$\hat{y} = x_0^T \beta + \sum_{i=1}^N \ell_i(x_0) \varepsilon_i,$$

where  $\ell_i(x_0)$  is the  $i^{\text{th}}$  element of  $X(X^T X)^{-1} x_0$ .



$$Y = X^T \beta + \varepsilon,$$

$$\hat{y} = x_0^T \beta + \sum_{i=1}^N \ell_i(x_0) \varepsilon_i,$$

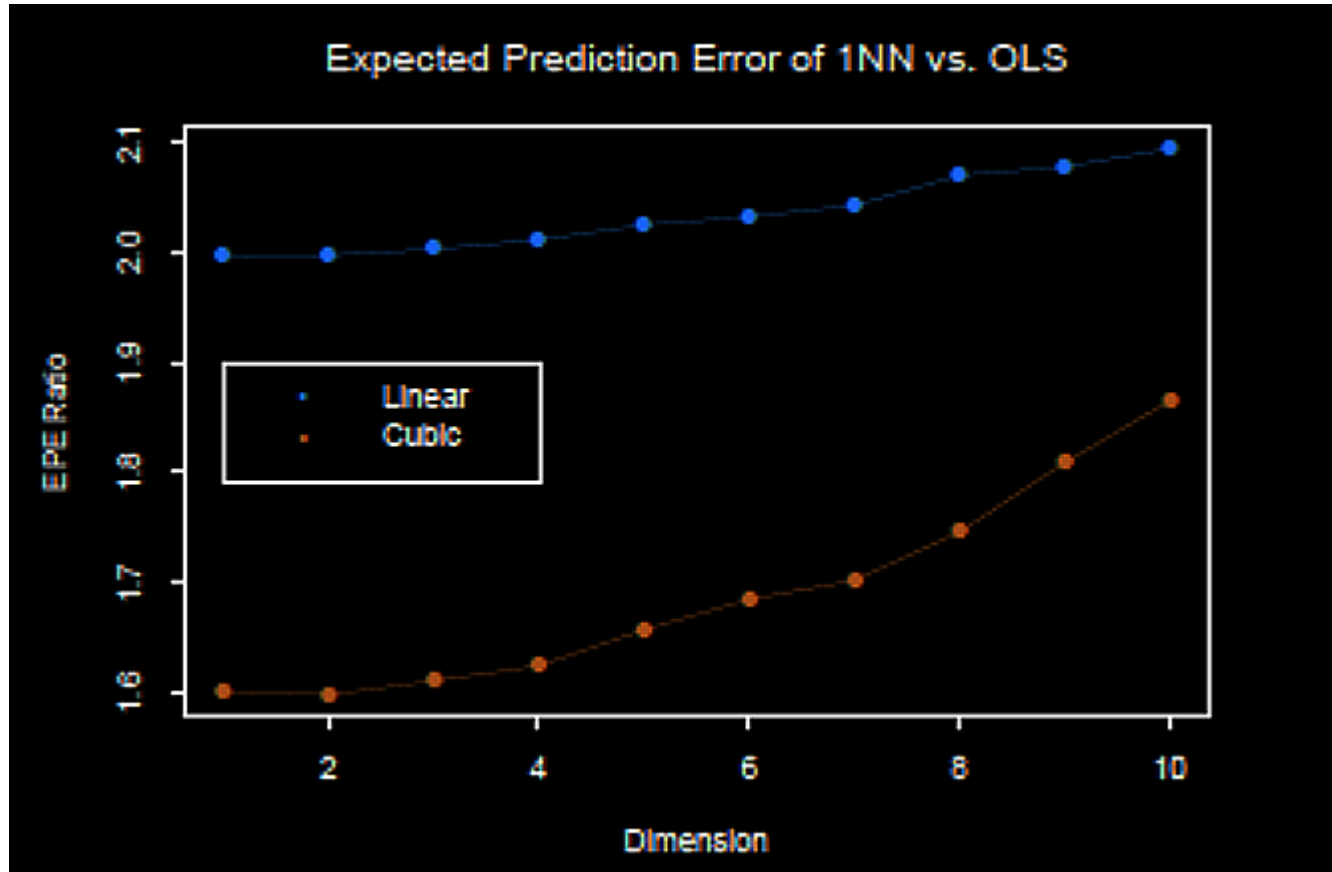
$$\begin{aligned} \text{EPE}(x_0) &= \mathbb{E}_{y_0|x_0} \mathbb{E}_{\mathcal{T}} (y_0 - \hat{y}_0)^2 \\ &= \text{Var}(y_0|x_0) + \mathbb{E}_{\mathcal{T}} [\hat{y}_0 - \mathbb{E}_{\mathcal{T}} \hat{y}_0]^2 + [\mathbb{E}_{\mathcal{T}} \hat{y}_0 - x_0^T \beta]^2 \\ &= \text{Var}(y_0|x_0) + \text{Var}_{\mathcal{T}}(\hat{y}_0) + \text{Bias}^2(\hat{y}_0) \\ &= \sigma^2 + \mathbb{E}_{\mathcal{T}} x_0^T (X^T X)^{-1} x_0 \sigma^2 + 0^2. \end{aligned}$$

If  $N$  is large and  $T$  were selected at random, and assuming  $E(X) = 0$ ,  
then  $X^T X \rightarrow N \text{Cov}(X)$  and

$$\begin{aligned} E_{x_0} EPE(x_0) &\sim E_{x_0} x_0^T \text{Cov}(X)^{-1} x_0 \sigma^2 / N + \sigma^2 \\ &= \text{trace}[\text{Cov}(X)^{-1} \text{Cov}(x_0)] \sigma^2 / N + \sigma^2 \\ &= \sigma^2 (p/N) + \sigma^2. \end{aligned} \tag{2.28}$$

- we see that the expected  $EPE$  increases linearly as a function of  $p$ , with slope  $\sigma^2/N$ . If  $N$  is large and/or  $\sigma^2$  is small, this growth in variance is negligible (0 in the deterministic case). By imposing some heavy restrictions on the class of models being fitted, we have avoided the curse of dimensionality.

# Expected Prediction Error of 1NN vs. OLS



- **FIGURE 2.9.** The curves show the expected prediction error (at  $x_0 = 0$ ) for 1-nearest neighbor relative to least squares for the model  $Y = f(X) + \varepsilon$ . For the orange curve,  $f(x) = x_1$ , while for the blue curve  $f(x) = \frac{1}{2}(x_1 + 1)^3$ .

- **Figure 2.9** compares 1 –nearest neighbor vs. Least squares in two situations, both of which have the form  $Y = f(X) + \varepsilon, X$  uniform as before, and  $\varepsilon \sim N(0, 1)$ . The sample size is  $N = 500$ .
- For the orange curve,  $f(x)$  is linear in the first coordinate, for the blue curve, cubic as in Figure 2.8. Shown is the relative *EPE* of 1-nearest neighbor to least squares, which appears to start at around 2 for the linear case. Least squares is unbiased in this case, and as discussed above the *EPE* is slightly above  $\sigma^2 = 1$ .
- The *EPE* for 1-nearest neighbor is always above 2, since the variance of  $\hat{f}(x_0)$  in this case is at least  $\sigma^2$ , and the ratio increases with dimension as the nearest neighbor strays from the target point. For the cubic case, least squares is biased, which moderates the ratio. Clearly we could manufacture examples where the bias of least squares would dominate the variance, and the 1-nearest neighbor would come out the winner.

# Statistical Models, Supervised Learning and Function Approximation

- Squared error loss lead us to the regression function  $f(x) = E(Y | X = x)$  for a quantitative response. The class of nearest-neighbor methods can be viewed as direct estimates of this conditional expectation, but we have seen that they can fail in at least two ways:
  - ❖ If the dimension of the input space is high, the nearest neighbors need not be close to the target point, and can result in large errors;
  - ❖ If special structure is known to exist, this can be used to reduce both the bias and the variance of the estimates.

# A Statistical Model for the Joint Distribution $Pr(X, Y)$

- Suppose in fact that our data arose from a statistical model.

$$Y = f(X) + \varepsilon,$$

- where the random error  $\varepsilon$  has  $E(\varepsilon) = 0$  and is independent of  $X$ . Note that for this model,  $f(x) = E(Y | X = x)$ , and in fact the conditional distribution  $Pr(Y | X)$  depends on  $X$  only through the conditional mean  $f(x)$ .
- The additive error model is a useful approximation to the truth. For most systems the input–output pairs  $(X, Y)$  will not have a deterministic relationship  $Y = f(X)$ . Generally there will be other unmeasured variables that also contribute to  $Y$ , including measurement error. The additive model assumes that we can capture all these departures from a deterministic relationship via the error  $\varepsilon$ .

# Supervised Learning

- Assembles a training set of observations  $T = (x_i, y_i), i = 1, \dots, N$ . The observed input values to the system  $x_i$  are also fed into an artificial system, known as a learning algorithm (usually a computer program), which also produces outputs  $\hat{f}(x_i)$  in response to the inputs. The learning algorithm has the property that it can modify its input/output relationship  $\hat{f}$  in response to differences  $y_i - \hat{f}(x_i)$  between the original and generated outputs. This process is known as learning by example. Upon completion of the learning process the hope is that the artificial and real outputs will be close enough to be useful for all sets of inputs likely to be encountered in practice.

# Function Approximation

- The approach taken in applied mathematics and statistics has been from the perspective of function approximation and estimation. Here the data pairs  $\{x_i, y_i\}$  are viewed as points in a  $(p + 1)$ -dimensional Euclidean space. The function  $f(x)$  has domain equal to the  $p$ -dimensional input subspace, and is related to the data via a model  $y_i = f(x_i) + \varepsilon_i$ . Although somewhat less glamorous than the learning paradigm, treating supervised learning as a problem in function approximation encourages the geometrical concepts of Euclidean spaces and mathematical concepts of probabilistic inference to be applied to the problem.
- Many of the approximations we will encounter have associated a set of parameters  $\theta$  that can be modified to suit the data at hand. For example, the linear model  $f(x) = x^T \beta$  has  $\theta = \beta$ . Another class of useful approximates can be expressed as *linear basis expansions*.

$$f_{\theta}(x) = \sum_{k=1}^K h_k(x) \theta_k,$$



- where the  $h_k$  are a suitable set of functions or transformations of the input vector  $x$ . Traditional examples are polynomial and trigonometric expansions, where for example  $h_k$  might be  $x_1^2$ ,  $x_1x_2^2$ ,  $\cos(x_1)$  and so on. Encounter nonlinear expansions, such as the sigmoid transformation common to neural network models,

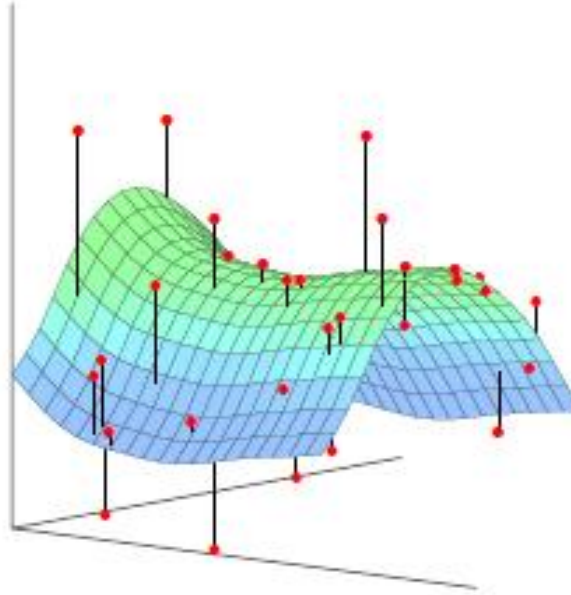
$$h_k(x) = \frac{1}{1 + \exp(-x^T \beta_k)}$$

- We can use least squares to estimate the parameters  $\theta$  in  $f_\theta$  as we did for the linear model, by minimizing the residual sum-of-squares

$$RSS(\theta) = \sum_{i=1}^N (y_i - f_\theta(x_i))^2$$

- As a function of  $\theta$  This seems a reasonable criterion for an additive error model. In terms of function approximation, we imagine our parameterized function as a surface in  $p + 1$  space, and what we observe are noisy realizations from it. This is easy to visualize when  $p = 2$  and the vertical coordinate is the output  $y$ , as in Figure 2.10. The noise is in the output coordinate, so we find the set of parameters such that the fitted surface gets as close to the observed points as possible, where close is measured by the sum of squared vertical errors in  $RSS(\theta)$ .
- For the linear model we get a simple closed form solution to the minimization problem. This is also true for the basis function methods, if the basis functions themselves do not have any hidden parameters. Otherwise the solution requires either iterative methods or numerical optimization.
- While least squares is generally very convenient, it is not the only criterion used and in some cases would not make much sense.

# Statistical Models, Supervised Learning and Function Approximation



- **FIGURE 2.10.** *Least squares fitting of a function of two inputs. The parameters of  $f_{\theta}(x)$  are chosen so as to minimize the sum-of-squared vertical errors.*

- Principle for estimation is maximum likelihood estimation. Suppose we have a random sample  $y_i, i = 1, \dots, N$  from a density  $\text{Pr}(y)$  indexed by some parameters  $\theta$ . The log-probability of the observed sample is

$$L(\theta) = \sum_{i=1}^N \log \text{Pr}_{\theta}(y_i).$$

- The principle of maximum likelihood assumes that the most reasonable values for  $\theta$  are those for which the probability of the observed sample is largest. Least squares for the additive error model  $Y = f_{\theta}(X) + \varepsilon$ , with  $\varepsilon \sim N(0, \sigma^2)$ , is equivalent to maximum likelihood using the conditional likelihood

$$\text{Pr}(Y | X, \theta) = N(f_{\theta}(X), \sigma^2).$$

- So although the additional assumption of normality seems more restrictive, the results are the same. The log-likelihood of the data is

$$L(\theta) = -\frac{N}{2} \log(2\pi) - N \log \sigma - \frac{1}{2\sigma^2} \sum_{i=1}^N (y_i - f_{\theta}(x_i))^2$$

- And the only term involving  $\theta$  is the last, which is  $RSS(\theta)$  up to a scalar negative multiplier.
- A more interesting example is the multinomial likelihood for the regression function  $\Pr(G|X)$  for a qualitative output  $G$ . Suppose we have a model  $\Pr(G = G_k|X = x) = p_{k,\theta}(x), k = 1, \dots, K$  for the conditional probability of each class given  $X$ , indexed by the parameter vector  $\theta$ .
- Likelihood (also referred to as the cross-entropy) is

$$L(\theta) = \sum_{i=1}^N \log p_{g_i, \theta}(x_i)$$

# Structured Regression Models

- We introduces classes of such structured approaches.

## ❖ Difficulty of the Problem

- Consider the RSS criterion for an arbitrary function  $f$ ,

$$RSS(f) = \sum_{i=1}^N (y_i - f(x_i))^2$$

- Minimizing (2.37) leads to infinitely many solutions: any function  $\hat{f}$  passing through the training points  $(x_i, y_i)$  is a solution. Any particular solution chosen might be a poor predictor at test points different from the training points.
- In order to obtain useful results for finite  $N$ , we must restrict the eligible solutions to (2.37) to a smaller set of functions. How to decide on the nature of the restrictions is based on considerations outside of the data. These restrictions are sometimes encoded via the parametric representation of  $f_\theta$ , or may be built into the learning method itself, either implicitly or explicitly.

- Any restrictions imposed on  $f$  that lead to a unique solution to 2.38 do not really remove the ambiguity 2.8 Classes of Restricted Estimators 33 caused by the multiplicity of solutions. There are infinitely many possible restrictions, each leading to a unique solution, so the ambiguity has simply been transferred to the choice of constraint.
- In general the constraints imposed by most learning methods can be described as complexity restrictions of one kind or another. This usually means some kind of regular behavior in small neighborhoods of the input space.

- The strength of the constraint is dictated by the neighborhood size. The larger the size of the neighborhood, the stronger the constraint, and the more sensitive the solution is to the particular choice of constraint. For example, local constant fits in infinitesimally small neighborhoods is no constraint at all; local linear fits in very large neighborhoods is almost a globally linear model, and is very restrictive.
- The nature of the constraint depends on the metric used. Some methods, such as kernel and local regression and tree-based methods, directly specify the metric and size of the neighborhood. The nearest-neighbor methods discussed so far are based on the assumption that locally the function is constant; close to a target input  $x_0$ , the function does not change much, and so close outputs can be averaged to produce  $\hat{f}(x_0)$ . Other methods such as splines, neural networks and basis-function methods implicitly define neighborhoods of local behavior. In Section 5.4.1 we discuss the concept of an equivalent kernel.



- Attempts to produce locally varying functions in small isotropic neighborhoods will run into problems in high dimensions—again the curse of dimensionality. And conversely, all methods that overcome the dimensionality problems have an associated—and often implicit or adaptive—metric for measuring neighborhoods, which basically does not allow the neighborhood to be simultaneously small in all directions.

## **Classes of Restricted Estimators**

- The variety of nonparametric regression techniques or learning methods fall into a number of different classes depending on the nature of the restrictions imposed. Brief summary, detailed descriptions are given in later chapters. *Smoothing* parameters, that control the effective size of the local neighborhood. Three broad classes.

# Roughness Penalty and Bayesian Methods

- Here the class of functions is controlled by explicitly penalizing  $RSS(f)$  with a roughness penalty.

$$PRSS(f; \lambda) = RSS(f) + \lambda J(f).$$

## Kernel Methods and Local Regression

- These methods can be thought of as explicitly providing.

# Basis Functions and Dictionary Methods

- Radial basis functions are symmetric p-dimensional kernels located at particular centroids,

- $$f_{\theta}(x) = K_{\lambda_m}(\mu_m, x)\theta_m;$$

- For example, the Gaussian kernel  $K_{\lambda}(\mu, x) = e^{-\frac{\|x-\mu\|^2}{2\lambda}}$  is popular.
- Radial basis functions have centroids  $\mu_m$  and scales  $\lambda_m$  that have to be determined. The spline basis functions have knots.
- A single-layer feed-forward neural network model with linear output weights can be thought of as an adaptive basis function method. The model has the form

$$f_{\theta}(x) = \sum_{m=1}^M \beta_m \sigma(\alpha_m^T x + b_m),$$

# Model Selection and the Bias–Variance Tradeoff

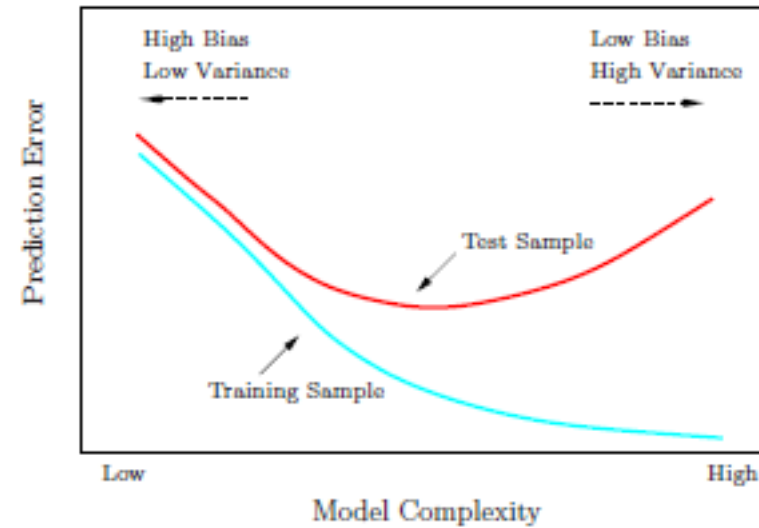
- All the models described above and many others discussed in later chapters have a smoothing or complexity parameter that has to be determined:
  - ❖ the multiplier of the penalty term;
  - ❖ the width of the kernel;
  - ❖ or the number of basis functions.
- The data arise from a model  $Y = f(X) + \varepsilon$ , with  $E(\varepsilon) = 0$  and
- $Var(\varepsilon) = \sigma^2$ . The expected prediction error at  $x_0$ , also known as test or generalization error, can be decomposed:

$$\begin{aligned} EPE_k(x_0) &= E[(Y - \hat{f}_k(x_0))^2 | X = x_0] \\ &= \sigma^2 + [Bias^2(\hat{f}_k(x_0)) + VarT(\hat{f}_k(x_0))] \end{aligned} \quad 2.46$$

$$= \sigma^2 + \left[ f(x_0) - \frac{1}{k} \sum_{\ell=1}^k f(x_{(\ell)}) \right]^2 + \frac{\sigma^2}{k}. \quad 2.47$$

- The first term  $\sigma^2$  is the *irreducible Error*.
- The second and third terms are under our control, and make up the mean squared error of  $\hat{f}_k(x_0)$  in estimating  $f(x_0)$ , which is broken down into a bias component and a variance component. The bias term is the squared difference between the true mean  $f(x_0)$  and the expected value of the estimate— $\left[ E_T \left( \hat{f}_k(x_0) \right) - f(x_0) \right]^2$ —where the expectation averages the randomness in the training data. This term will most likely increase with  $k$ , if the true function is reasonably smooth. For small  $k$  the few closest neighbors will have values  $f(x_{(\ell)})$  close to  $f(x_0)$ , so their average should be close to  $f(x_0)$ . As  $k$  grows, the neighbors are further away, and then anything can happen.

# Model Selection and the Bias–Variance Tradeoff



**FIGURE 2.11.** Test and training error as a function of model complexity.

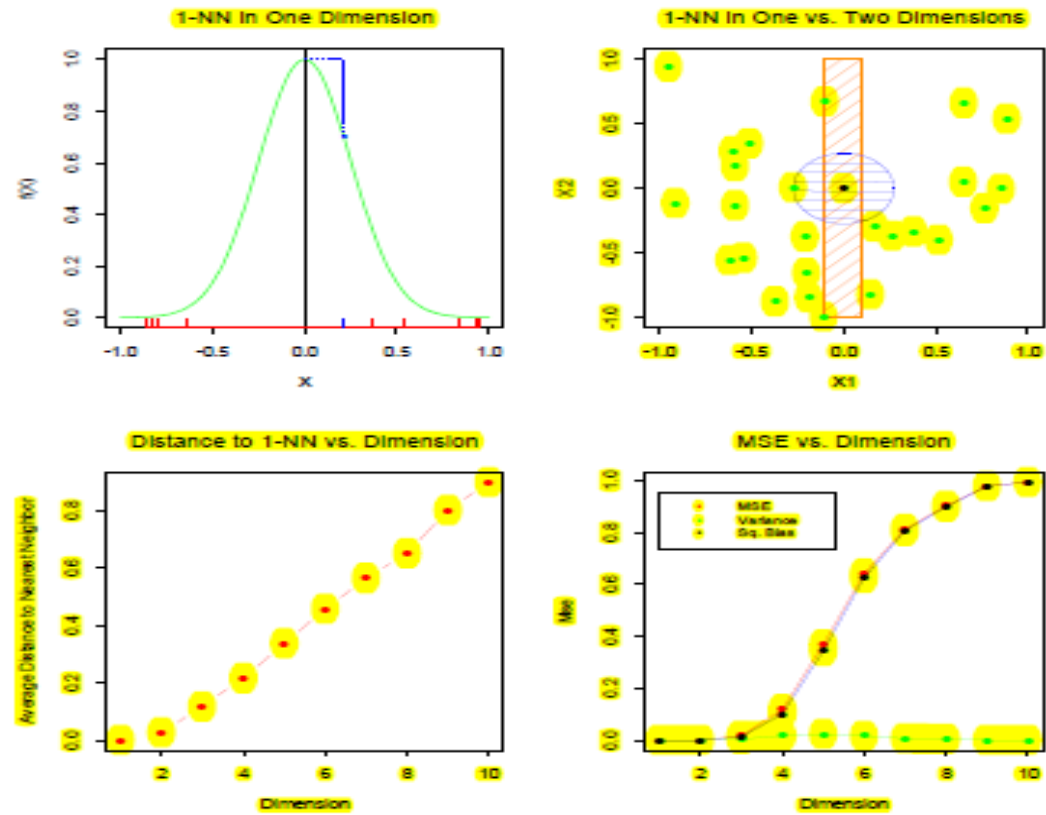
- **Figure 2.11** shows the typical behavior of the test and training error, as model complexity is varied. The training error tends to decrease whenever we increase the model complexity, that is, whenever we fit the data harder. However with too much fitting, the model adapts itself too closely to the training data, and will not generalize well (i.e., have large test error). In that case the predictions  $\hat{f}(x_0)$  will have large variance, as reflected in the last term of expression (2.46). In contrast, if the model is not complex enough, it will under fit and may have large bias, again resulting in poor generalization.

- The variance term is simply the variance of an average here, and decreases as the inverse of  $k$ . So as  $k$  varies, there is a bias–variance tradeoff.
- More generally, as the model complexity of our procedure is increased, the variance tends to increase and the squared bias tends to decrease. The opposite behavior occurs as the model complexity is decreased. For  $k$ -nearest neighbors, the model complexity is controlled by  $k$ .
- Typically we would like to choose our model complexity to trade bias off with variance in such a way as to minimize the test error. An obvious estimate of test error is the training error  $\frac{1}{N} \sum_i (y_i - \hat{y}_i)^2$ . Unfortunately training error is not a good estimate of test error, as it does not properly account for model complexity.





# Local Methods in High Dimensions



**FIGURE 2.7.** A simulation example, demonstrating the curse of dimensionality and its effect on MSE, bias and variance.

# Local Methods in High Dimensions