

27 Mean-Square-Error Inference

Inference deals with the estimation of hidden parameters or random variables from observations of other related variables. In this chapter, we study the basic, yet fundamental, problem of inferring an unknown random quantity from observations of another random quantity by using the *mean-square-error* (MSE) criterion. Several other design criteria can be used for inference purposes besides MSE, such as the mean-absolute error (MAE) and the maximum a-posteriori (MAP) criteria. We will encounter these possibilities in future chapters, starting with the next chapter. We initiate our discussions of inference problems though by focusing on the MSE criterion due to its mathematical tractability and because it sheds light on several important questions that arise in the study of inference problems in general.

In our treatment of inference problems, we will encounter three broad formulations:

- (a) In one instance, we will model the unknown as a *random* variable, denoted by the boldface symbol \mathbf{x} . The objective will be to predict (or infer) the value of \mathbf{x} from observations of a related variable \mathbf{y} . The predictor or estimator for \mathbf{x} will be denoted by $\hat{\mathbf{x}}$. Different design criteria will lead to different constructions for $\hat{\mathbf{x}}$. In this chapter, we discuss the popular MSE criterion.
- (b) In another instance, we will continue to model the unknown \mathbf{x} as a random variable but will limit its values to *discrete* levels, such as having \mathbf{x} assume the values $+1$ or -1 . This type of formulation is prevalent in classification problems and will be discussed at length in future chapters, starting with the next chapter, where we examine the Bayes classifier.
- (c) In a third instance, we will model the unobservable as an unknown *constant*, denoted by the Greek symbol θ , rather than a random variable. This type of problem is frequent in applications requiring fitting models onto data and will be discussed in future chapters, for example, when we introduce the maximum-likelihood and expectation maximization paradigms.

27.1 INFERENCE WITHOUT OBSERVATIONS

We consider first a simple yet useful estimation problem, which relates to estimating a random variable $\mathbf{x} \in \mathbb{R}$ when the only information that is available about \mathbf{x} is its mean, \bar{x} . Our objective is to estimate the value that \mathbf{x} will assume in a given experiment. We denote the *estimate* for \mathbf{x} by the notation \hat{x} ; it is a *deterministic* quantity (i.e., a number). But how do we come up with a value for \hat{x} ? And how do we decide whether this value is optimal or not? And if optimal, in what sense? These inquiries are at the heart of every inference problem. To answer them, we first need to choose a cost (also called *risk*) function to penalize the estimation error. The resulting estimate \hat{x} will be optimal only in the sense that it leads to the smallest cost or risk value. Different choices for the cost function will generally lead to different choices for \hat{x} , each of which will be optimal in its own way.

27.1.1 Problem Formulation

The design criterion we study first is the famed MSE criterion; several other criteria are possible and we discuss other important choices in the next chapter. The MSE criterion is based on introducing the error signal:

$$\tilde{\mathbf{x}} \triangleq \mathbf{x} - \hat{x} \quad (27.1)$$

and on determining \hat{x} by minimizing the MSE, which is defined as the mean of the squared error $\tilde{\mathbf{x}}^2$:

$$\hat{x} \triangleq \underset{\hat{x}}{\operatorname{argmin}} \mathbb{E} \tilde{\mathbf{x}}^2 \quad (27.2)$$

The error $\tilde{\mathbf{x}}$ is a random variable since \mathbf{x} is random. The resulting estimate, \hat{x} , is called the *least-mean-squares estimate* (LMSE) of \mathbf{x} . For added emphasis, we could have written \hat{x}_{MSE} , with the subscript MSE, in order to highlight the fact that this is an estimate for \mathbf{x} that is based on minimizing the MSE criterion defined by (27.2). This notation is unnecessary in this chapter because we will be discussing mainly the MSE criterion. However, in later chapters, when we introduce other design criteria, it will become necessary to use the MSE subscript to distinguish the MSE estimate, \hat{x}_{MSE} , from other estimates such as \hat{x}_{MAP} , \hat{x}_{MAE} , or \hat{x}_{ML} , where the subscripts MAP, MAE, and ML will be referring to maximum a-posteriori, mean-absolute error, and maximum-likelihood estimators.

Returning to (27.2), it is immediate to verify that the solution \hat{x} is given by

$$\hat{x} = \bar{x} \quad (27.3)$$

and that the resulting minimum mean-square error (MMSE) is

$$\boxed{\text{MMSE} = \mathbb{E} \tilde{\mathbf{x}}^2 = \sigma_x^2} \quad (27.4)$$

We therefore say that the best estimate of a random variable (best from the perspective of MSE estimation) when only its mean is known is the mean value of the random variable itself.

Proof of (27.3)–(27.4) using differentiation: We expand the MSE from (27.2) to get

$$\mathbb{E} \tilde{\mathbf{x}}^2 = \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}})^2 = \mathbb{E} \mathbf{x}^2 - 2\bar{x}\hat{\mathbf{x}} + \hat{\mathbf{x}}^2 \quad (27.5)$$

and then differentiate the right-hand side with respect to the unknown $\hat{\mathbf{x}}$. Setting the derivative to zero gives

$$\frac{\partial \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}})^2}{\partial \hat{\mathbf{x}}} = -2\bar{x} + 2\hat{\mathbf{x}} = 0 \implies \hat{\mathbf{x}} = \bar{x} \quad (27.6)$$

This choice for $\hat{\mathbf{x}}$ minimizes the MSE since the risk (27.5) is quadratic in $\hat{\mathbf{x}}$. ■

Proof of (27.3)–(27.4) using completion-of-squares: An alternative argument to arrive at the same conclusion relies on the use of a completion-of-squares step. It is immediate to verify that the risk function in (27.2) can be rewritten in the following form, by adding and subtracting \bar{x} :

$$\begin{aligned} \mathbb{E} \tilde{\mathbf{x}}^2 &= \mathbb{E} ((\mathbf{x} - \bar{x}) + (\bar{x} - \hat{\mathbf{x}}))^2 \\ &= \mathbb{E} (\mathbf{x} - \bar{x})^2 + (\bar{x} - \hat{\mathbf{x}})^2 + 2 \underbrace{\mathbb{E} (\mathbf{x} - \bar{x})(\bar{x} - \hat{\mathbf{x}})}_{=0} \\ &= \sigma_x^2 + (\bar{x} - \hat{\mathbf{x}})^2 \end{aligned} \quad (27.7)$$

The above result expresses the risk as the sum of two *nonnegative* terms where only the second term depends on the unknown, $\hat{\mathbf{x}}$. It is clear that the choice $\hat{\mathbf{x}} = \bar{x}$ annihilates the second term and results in the smallest possible value for the MSE. This value is referred to as the MMSE and is equal to σ_x^2 . ■

27.1.2 Interpretation

There are good reasons for using the MSE criterion (27.2). The simplest one perhaps is that the criterion is amenable to mathematical manipulations, more so than any other criterion. In addition, the criterion is attempting to force the estimation error, $\tilde{\mathbf{x}}$, to assume values close to its mean, which is zero since

$$\mathbb{E} \tilde{\mathbf{x}} = \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}}) = \mathbb{E} (\mathbf{x} - \bar{x}) = \bar{x} - \bar{x} = 0 \quad (27.8)$$

Therefore, by minimizing $\mathbb{E} \tilde{\mathbf{x}}^2$, we are in effect minimizing the variance of the error, $\tilde{\mathbf{x}}$. Then, in view of the discussion in Section 3.2 regarding the interpretation of the variance of a random variable, we find that the MSE criterion is attempting to increase the likelihood of small errors.

The effectiveness of the estimation procedure (27.2) can be measured by examining the value of the resulting minimum risk, which is the variance of the estimation error, denoted by $\sigma_{\hat{x}}^2 = \mathbb{E} \tilde{\mathbf{x}}^2$. The above discussion tells us that

$$\sigma_{\hat{x}}^2 = \sigma_x^2 \quad (27.9)$$

so that the estimate $\hat{x} = \bar{x}$ *does not* reduce our initial uncertainty about \mathbf{x} : The error variable has the same variance as \mathbf{x} itself. We therefore find that, in this initial scenario, the performance of the MSE design procedure is limited. We are interested in estimation procedures that result in error variances that are smaller than the original signal variance. We discuss one such procedure in the next section.

The reason for the poor performance of the estimate $\hat{x} = \bar{x}$ lies in the lack of more sophisticated prior information about \mathbf{x} . Note that result (27.3) simply tells us that the best we can do, in the absence of any other information about a random variable \mathbf{x} , other than its mean, is to use the mean value of \mathbf{x} as our estimate. This result is at least intuitive. After all, the mean value of a random variable is, by definition, an indication of the value that we expect to observe on average in repeated experiments. Hence, in answer to the question, “*What is the best guess for \mathbf{x} ?*” the analysis tells us that the best guess is what we would expect for \mathbf{x} on average! This is a circular answer, but one that is at least consistent with intuition.

Example 27.1 (Guessing the class of an image) Assume a box contains an equal number of images of cats and dogs. An image is selected at random from the box and a random variable \mathbf{x} is associated with this experiment. The variable \mathbf{x} assumes the value $x = +1$ if the selected image is a cat and it assumes the value $x = -1$ otherwise. We say that \mathbf{x} represents the class (or label) variable: Its value specifies whether the selected image belongs to one class (+1 corresponding to cats) or the other (−1 corresponding to dogs). It is clear that \mathbf{x} is a *binary* random variable assuming the values ± 1 with equal probability. Then,

$$\bar{x} = \frac{1}{2} \times (+1) + \frac{1}{2} \times (-1) = 0 \quad (27.10)$$

$$\sigma_x^2 = \mathbb{E} \mathbf{x}^2 = 1 \quad (27.11)$$

If a user were to predict the class of the image that will be selected from the box beforehand then, according to the MSE criterion, the best estimate for \mathbf{x} will be $\hat{x} = \bar{x} = 0$. This estimate value is neither +1 nor −1. This example shows that the LMSE estimate does not always lead to a meaningful solution! In this case, using $\hat{x} = 0$ is not useful in guessing whether the realization for \mathbf{x} will be +1 or −1. If we could incorporate into the design procedure some additional information, besides the mean of \mathbf{x} , then we could perhaps come up with a better prediction for the class of the image.

Example 27.2 (Predicting a crime statistic) The US Federal Bureau of Investigation (FBI) publishes statistics on the crime rates in the country on an annual basis. Figure 27.1 plots the burglary rates per 100,000 inhabitants for the period 1997–2016. Assume we model the annual burglary rate as a random variable \mathbf{x} with some mean \bar{x} . By examining the plot, we find that this assumption is more or less reasonable only over the shorter range 2000–2009 during which the burglary rate remained practically flat

with fluctuations around some nominal average value. The rates are declining before 2000 and after 2010. Assume we did not know the burglary rate for the year 2010 and wanted to predict its value from the burglary rates observed in prior years. In this example, the probability distribution of \mathbf{x} is not known to evaluate its mean \bar{x} . Instead, we have access to measurements for the years 1997–2015. We can use the data from the years 2000–2009 to compute a sample mean and use it to predict $x(2010)$:

$$\hat{x}(2010) = \frac{1}{10} \sum_{n=2000}^{2009} x(n) \approx 732.6 \quad (27.12)$$

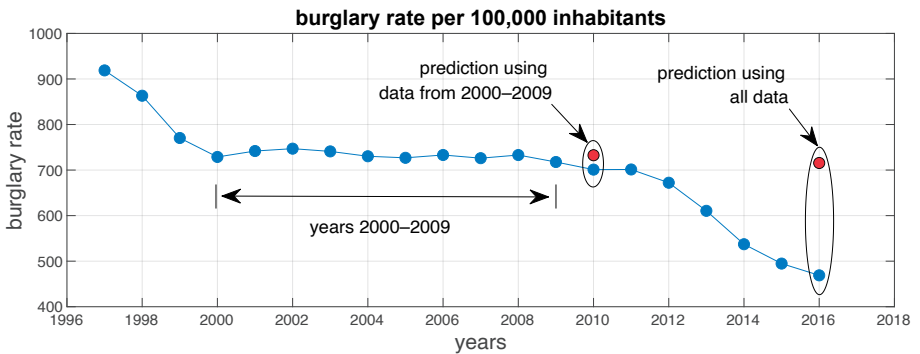


Figure 27.1 Plot of the annual burglary rates per 100,000 inhabitants in the United States during 1997–2016. Source: US Criminal Justice Information Services Division. Data found at <https://ucr.fbi.gov/crime-in-the-u.s/2016/crime-in-the-u.s.-2016/tables/table-1>.

The value 732.6 is close enough to the actual burglary rate observed for 2010, which is 701. If we were instead to predict the burglary rate for the year 2016 by using the data for the entire period 1997–2015, we would end up with

$$\hat{x}(2016) \approx 715.5 \quad (27.13)$$

which is clearly a bad estimate since the actual value is 468.9.

This example illustrates the fact that we will often be dealing with distributions that vary (i.e., drift) over time for various reasons, such as changing environmental conditions or, in the case of this example, crime deterrence policies that may have been put in place. This possibility necessitates the development of inference techniques that are able to adapt to variations in the statistical properties of the data in an *automated* manner. In this example, the statistical properties of the data during the period 2000–2009 are clearly different from the periods before 2000 and after 2009.

27.2 INFERENCE WITH OBSERVATIONS

Let us examine next the case in which more is known about the random variable \mathbf{x} , beyond its mean. Let us assume that we have access to an observation of a second random variable \mathbf{y} that is related to \mathbf{x} in some way. For example, \mathbf{y} could

be a noisy measurement of \mathbf{x} , say, $\mathbf{y} = \mathbf{x} + \mathbf{v}$, where \mathbf{v} denotes the disturbance, or \mathbf{y} could be the sign of \mathbf{x} , or dependent on \mathbf{x} in some other way.

27.2.1 Conditional Mean Estimator

Given two dependent random variables $\{\mathbf{x}, \mathbf{y}\}$, we pose the problem of determining the LMSE estimator for \mathbf{x} from an observation of \mathbf{y} . Observe that we are now employing the terminology “estimator” of \mathbf{x} as opposed to “estimate” of \mathbf{x} . In order to highlight this distinction, we denote the estimator of \mathbf{x} by the boldface notation $\hat{\mathbf{x}}$; it is now a random variable since it will be a function of \mathbf{y} , i.e.,

$$\hat{\mathbf{x}} = c(\mathbf{y}) \tag{27.14}$$

for some function $c(\cdot)$ to be determined. Once the function $c(\cdot)$ has been determined, evaluating it at a particular observation for \mathbf{y} , say, at $\mathbf{y} = y$, will result in an estimate for \mathbf{x} , i.e.,

$$\hat{\mathbf{x}} = c(\mathbf{y}) \Big|_{\mathbf{y}=y} = c(y) \tag{27.15}$$

Different occurrences for \mathbf{y} will lead to different estimates $\hat{\mathbf{x}}$. In Section 27.1 we did not need to make this distinction between an estimator $\hat{\mathbf{x}}$ and an estimate $\hat{\mathbf{x}}$. There we sought directly an estimate $\hat{\mathbf{x}}$ for \mathbf{x} since we did not have access to a random variable \mathbf{y} ; we only had access to the deterministic quantity \bar{x} .

The criterion we use to determine the estimator $\hat{\mathbf{x}}$ will continue to be the same MSE criterion. We define the error signal:

$$\tilde{\mathbf{x}} \triangleq \mathbf{x} - \hat{\mathbf{x}} \tag{27.16}$$

and determine $\hat{\mathbf{x}}$ by minimizing the MSE over all possible choices for the function $c(\cdot)$:

$$\min_{c(\cdot)} \mathbb{E} \tilde{\mathbf{x}}^2, \quad \text{subject to } \hat{\mathbf{x}} = c(\mathbf{y}) \tag{27.17}$$

We establish in the following that the solution of (27.17) is given by the conditional mean estimator:

$$\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x} | \mathbf{y} = y) \tag{27.18}$$

That is, the optimal choice for $c(y)$ in (27.15) is

$$c^o(y) = \mathbb{E}(\mathbf{x} | \mathbf{y} = y) \tag{27.19}$$

Recall from (3.100) that for continuous random variables \mathbf{x} and \mathbf{y} , the conditional expectation is computed via the integral expression:

$$\mathbb{E}(\mathbf{x} | \mathbf{y} = y) = \int_{x \in \mathcal{X}} x f_{\mathbf{x} | \mathbf{y}}(x | y) dx \tag{27.20a}$$

over the domain of \mathbf{x} , while for discrete random variables:

$$\mathbb{E}(\mathbf{x}|\mathbf{y} = y) = \sum_x x \mathbb{P}(\mathbf{x} = x|\mathbf{y} = y) \tag{27.20b}$$

in terms of the conditional probability values and where the sum is over the possible realizations for \mathbf{x} . We continue our presentation by using the notation for continuous random variables without loss in generality.

Returning to (27.18), this estimator is obviously unbiased since from the result of Prob. 3.25 we know that

$$\mathbb{E} \hat{\mathbf{x}} = \mathbb{E} \left(\mathbb{E}(\mathbf{x}|\mathbf{y}) \right) = \mathbb{E} \mathbf{x} = \bar{\mathbf{x}} \tag{27.21}$$

Moreover, the resulting minimum cost or MMSE will be given by

$$\text{MMSE} \triangleq \mathbb{E} \tilde{\mathbf{x}}^2 = \sigma_x^2 - \sigma_x^2 \tag{27.22}$$

which is smaller than the earlier value (27.4). Result (27.18) states that the least-mean-squares estimator of \mathbf{x} is its conditional expectation given \mathbf{y} . This result is again intuitive. In answer to the question, “*What is the best guess for \mathbf{x} given that we observe \mathbf{y} ?*” the analysis tells us that the best guess is what we would expect for \mathbf{x} given the occurrence of \mathbf{y} !

Derivation of (27.18) using differentiation: Using again the result of Prob. 3.25 we have

$$\begin{aligned} \mathbb{E}(\mathbf{x} - \hat{\mathbf{x}})^2 &= \mathbb{E} \left\{ \mathbb{E} \left((\mathbf{x} - \hat{\mathbf{x}})^2 | \mathbf{y} \right) \right\} \\ &= \mathbb{E} \left\{ \mathbb{E} \left(\mathbf{x}^2 - 2\mathbf{x}\hat{\mathbf{x}} + \hat{\mathbf{x}}^2 | \mathbf{y} \right) \right\} \\ &= \mathbb{E} \left\{ \mathbb{E} \left(\mathbf{x}^2 - 2\mathbf{x}c(\mathbf{y}) + c^2(\mathbf{y}) | \mathbf{y} \right) \right\}, \text{ since } \hat{\mathbf{x}} = c(\mathbf{y}) \\ &= \mathbb{E} \left\{ \mathbb{E}(\mathbf{x}^2|\mathbf{y}) - 2c(\mathbf{y})\mathbb{E}(\mathbf{x}|\mathbf{y}) + c^2(\mathbf{y}) \right\} \end{aligned} \tag{27.23}$$

It is sufficient to minimize the inner expectation relative to $c(\mathbf{y})$, for any realization $\mathbf{y} = y$. Differentiating and setting the derivative to zero at the optimal solution $c^o(y)$ gives $c^o(y) = \mathbb{E}(\mathbf{x}|\mathbf{y} = y)$. ■

Derivation of (27.18) using completion-of-squares: In this second derivation, we will establish two useful intermediate results, namely, the orthogonality conditions (27.26) and (27.29). We again refer to the result of Prob. 3.25, which we write more explicitly as

$$\mathbb{E} \mathbf{x} = \mathbb{E}_{\mathbf{y}} \left(\mathbb{E}_{\mathbf{x}|\mathbf{y}}(\mathbf{x}|\mathbf{y}) \right) \tag{27.24}$$

where the outer expectation is relative to the distribution of \mathbf{y} , while the inner expectation is relative to the conditional distribution of \mathbf{x} given \mathbf{y} . It follows that, for any real-valued function of \mathbf{y} , say, $g(\mathbf{y})$,

$$\begin{aligned} \mathbb{E} \mathbf{x} g(\mathbf{y}) &= \mathbb{E}_{\mathbf{y}} \left\{ \mathbb{E}_{\mathbf{x}|\mathbf{y}} \left(\mathbf{x} g(\mathbf{y}) \mid \mathbf{y} \right) \right\} \\ &= \mathbb{E}_{\mathbf{y}} \left\{ \mathbb{E}_{\mathbf{x}|\mathbf{y}} \left(\mathbf{x} \mid \mathbf{y} \right) g(\mathbf{y}) \right\} \\ &= \mathbb{E}_{\mathbf{y}} \left\{ \mathbb{E}_{\mathbf{x}|\mathbf{y}} \left(\mathbf{x} \mid \mathbf{y} \right) \right\} g(\mathbf{y}) \end{aligned} \tag{27.25}$$

This means that, for any $g(\mathbf{y})$, it holds that

$$\mathbb{E} \left(\mathbf{x} - \mathbb{E}(\mathbf{x}|\mathbf{y}) \right) g(\mathbf{y}) = 0 \iff \tilde{\mathbf{x}} \perp g(\mathbf{y}) \quad \text{(orthogonality condition)} \tag{27.26}$$

Result (27.26) states that the error variable $\tilde{\mathbf{x}} = \mathbf{x} - \mathbb{E}(\mathbf{x}|\mathbf{y})$ is orthogonal to any function $g(\cdot)$ of \mathbf{y} ; we also represent this result by the compact notation $\tilde{\mathbf{x}} \perp g(\mathbf{y})$. However, since $\mathbf{x} - \mathbb{E}(\mathbf{x}|\mathbf{y})$ is zero mean, then it also holds that $\tilde{\mathbf{x}}$ is uncorrelated with $g(\mathbf{y})$.

Using this intermediate result, we return to the risk (27.17), add and subtract $\mathbb{E}(\mathbf{x}|\mathbf{y})$ to its argument, and express it as

$$\mathbb{E} (\mathbf{x} - \hat{\mathbf{x}})^2 = \mathbb{E} \left((\mathbf{x} - \mathbb{E}(\mathbf{x}|\mathbf{y}) + \mathbb{E}(\mathbf{x}|\mathbf{y}) - \hat{\mathbf{x}})^2 \right) \tag{27.27}$$

The term $\mathbb{E}(\mathbf{x}|\mathbf{y}) - \hat{\mathbf{x}}$ is a function of \mathbf{y} . Therefore, if we choose $g(\mathbf{y}) = \mathbb{E}(\mathbf{x}|\mathbf{y}) - \hat{\mathbf{x}}$, then from the orthogonality property (27.26) we conclude that

$$\mathbb{E} (\mathbf{x} - \hat{\mathbf{x}})^2 = \mathbb{E} \left(\mathbf{x} - \mathbb{E}(\mathbf{x}|\mathbf{y}) \right)^2 + \mathbb{E} \left(\mathbb{E}(\mathbf{x}|\mathbf{y}) - \hat{\mathbf{x}} \right)^2 \tag{27.28}$$

Now, only the second term on the right-hand side is dependent on $\hat{\mathbf{x}}$ and the MSE is minimized by choosing $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|\mathbf{y})$. To evaluate the resulting MMSE we first use the orthogonality property (27.26), along with the fact that $\hat{\mathbf{x}} = \mathbb{E}(\mathbf{x}|\mathbf{y})$ is itself a function of \mathbf{y} , to conclude that

$$\mathbb{E} (\mathbf{x} - \hat{\mathbf{x}}) \hat{\mathbf{x}} = 0 \iff \tilde{\mathbf{x}} \perp \hat{\mathbf{x}} \tag{27.29}$$

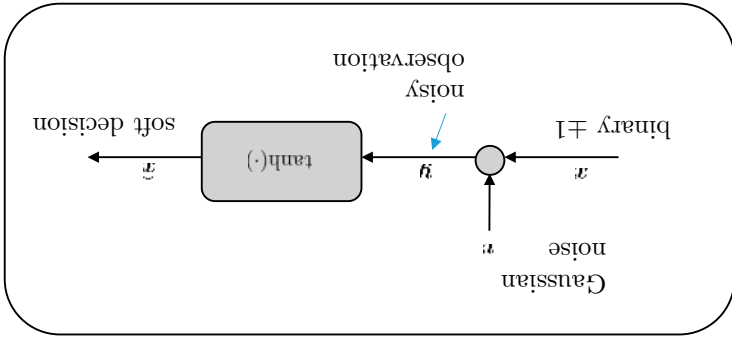
In other words, the estimation error, $\tilde{\mathbf{x}}$, is uncorrelated with the optimal estimator. Using this result, we can evaluate the MMSE as follows:

$$\begin{aligned} \mathbb{E} \tilde{\mathbf{x}}^2 &= \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}})(\mathbf{x} - \hat{\mathbf{x}}) \\ &= \mathbb{E} (\mathbf{x} - \hat{\mathbf{x}}) \mathbf{x} \quad \text{(because of (27.29))} \\ &= \mathbb{E} \mathbf{x}^2 - \mathbb{E} \hat{\mathbf{x}}(\tilde{\mathbf{x}} + \hat{\mathbf{x}}) \quad \text{(since } \mathbf{x} = \tilde{\mathbf{x}} + \hat{\mathbf{x}} \text{)} \\ &= \mathbb{E} \mathbf{x}^2 - \mathbb{E} \hat{\mathbf{x}}^2 \quad \text{(because of (27.29) again)} \\ &= (\mathbb{E} \mathbf{x}^2 - \bar{x}^2) + (\bar{x}^2 - \mathbb{E} \hat{\mathbf{x}}^2) \\ &= \sigma_x^2 - \sigma_{\hat{\mathbf{x}}}^2 \end{aligned} \tag{27.30}$$

■

Example 27.3 (From soft to hard decisions) Let us return to Example 27.1, where \mathbf{x} is a binary signal that assumes the values ± 1 with probability $1/2$. Recall that \mathbf{x} represents the class of the selected image (cats or dogs). In that example, we only assumed knowledge of the mean of \mathbf{x} and concluded that the resulting MSE estimate was not meaningful because it led to $\hat{\mathbf{x}} = 0$, which is neither $+1$ nor -1 . We are going to assume now that we have access to some additional information about \mathbf{x} . Specifically,

Figure 27.2 Estimation of a binary signal ± 1 observed under unit-variance additive Gaussian noise.



The result is represented in Fig. 27.2. If the variance of the measurement noise v were not fixed at 1 but denoted more generally by σ_v^2 , then the same argument would lead to $\hat{x} = \tanh(y/\sigma_v)$ with y scaled by σ_v^2 – see Prob. 27.2.

$$\hat{x} = \tanh(y) \triangleq \frac{e^y - e^{-y}}{e^y + e^{-y}} \tag{27.33}$$

Our intuition tells us that we should be able to do better here than in Example 27.1. But beware, even here, we will arrive at some interesting conclusions. According to (27.18), the optimal estimator for x given y is the conditional mean $\hat{x} = \mathbb{E}(x|y)$, which we evaluated earlier in (3.117) and determined that:

$$f^a(v) = \frac{1}{\sqrt{2\pi}} e^{-v^2/2} \tag{27.32}$$

We are then faced with the problem of deciding whether $x = \pm 1$ from the soft version y . Obviously, the nature of the perturbation v in (27.31) depends on the method that is used to generate the approximation y . In this example, and in order to keep the analysis tractable, we will assume that v and x are independent of each other and, moreover, that v has zero mean, unit variance, and is Gaussian-distributed with probability density function (pdf):

where the symbol v denotes the disturbance. How do measurements y of this type arise? We are going to encounter later in this text several inference techniques that provide approximate estimates for discrete variables. Rather than detect whether the unknown x is $+1$ or -1 (which we refer to as performing *hard* decisions), these other methods will return approximate values for x such as claiming that it is 0.9 or -0.7 (which we refer to as performing *soft* decisions). The soft value is a *real* (and not discrete) number and it can be interpreted as being a perturbed version of the actual label x .

$$y = x + v \tag{27.31}$$

we are going to assume that we have some noisy measurement of x , denoted by y , say, as

(27.33) tells the designer to estimate x by computing $\tanh(y)$. But, once again, this value will not be $+1$ or -1 ; it will be a real number somewhere inside the interval $(-1, 1)$. The designer will be induced to make a hard decision of the form:

$$\text{decide in favor of } \begin{cases} +1, & \text{if } \tanh(y) \text{ is nonnegative} \\ -1, & \text{if } \tanh(y) \text{ is negative} \end{cases} \quad (27.34)$$

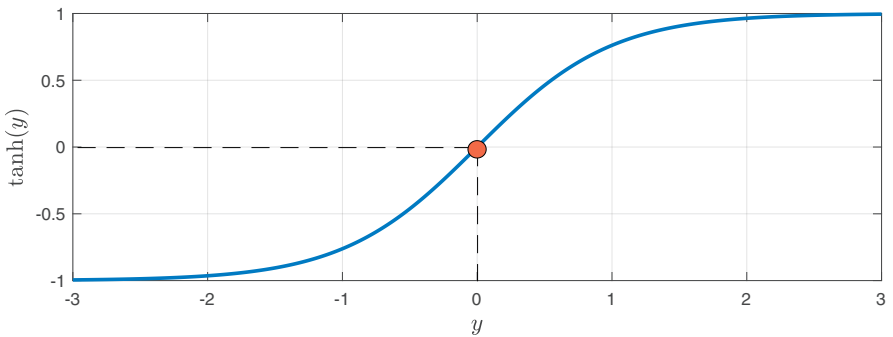


Figure 27.3 A plot of the hyperbolic tangent function, $\tanh(y)$. Observe that the curve tends to ± 1 as $y \rightarrow \pm\infty$.

In effect, the designer is implementing the alternative estimator:

$$\hat{x} = \text{sign}(\tanh(y)) \quad (27.35)$$

where $\text{sign}(\cdot)$ denotes the sign of its argument; it is equal to $+1$ if the argument is nonnegative and -1 otherwise:

$$\text{sign}(x) \triangleq \begin{cases} +1, & \text{if } x \geq 0 \\ -1, & \text{if } x < 0 \end{cases} \quad (27.36)$$

We therefore have a situation where the optimal estimator (27.33), although known in closed form, does not solve the original problem of recovering the symbols ± 1 directly. Instead, the designer is forced to implement a *suboptimal* solution; it is suboptimal from a LMSE point of view. Actually, the designer could consider implementing the following simpler suboptimal estimator directly:

$$\hat{x} = \text{sign}(y) \quad (27.37)$$

where the $\text{sign}(\cdot)$ function operates directly on y rather than on $\tanh(y)$ – see Fig. 27.4. Both suboptimal implementations (27.35) and (27.37) lead to the *same* result since, as is evident from Fig. 27.3:

$$\text{sign}(\tanh(y)) = \text{sign}(y) \quad (27.38)$$

We say that implementation (27.37) provides *hard decisions*, while implementation (27.33) provides *soft decisions*. We will revisit this problem later in Example 28.3 and show how the estimator (27.37) can be interpreted as being optimal relative to another design criterion; specifically, it will be the optimal Bayes classifier for the situation under study.

The purpose of Examples 27.1 and 27.3 is not to confuse the reader, but to stress the fact that an optimal estimator is optimal only in the sense that it satisfies a certain optimality criterion. One should not confuse an optimal guess with a perfect guess. One