

# Approximating unknown functions when fitting errors involve costs

Hans-Dieter Heike, Matei Demetrescu

**Abstract.** The problem of approximating unknown functions is discussed for the case where approximation errors are evaluated by means of a general cost-of-error (loss) function, not necessarily the squared-error one. To ensure minimal overall loss, approximation is carried out by fitting under the relevant loss function. Convergence results are provided in a general framework, allowing among others for Taylor approximations, approximations with Hermite polynomials or approximations with Neural Networks and for random measurement error.

**M.S.C. 2000:** 62G08, 65D10.

**Key words:** asymmetric loss, interpolation, sieve approximation.

## 1 Introduction

Function approximation is a common task in many branches of applied science, e.g. in control engineering or operations research. A specific response function is often desired, say  $f(x)$  with  $x \in D \subset \mathbb{R}$ , but, in many cases, the desired function is not available in closed form or is not easily handled. Also, it may be that the function is only known at some points (nodes)  $x_t$  of its support,  $t = 1, 2, \dots, T$ , and needs interpolation, resulting in the use of an approximating function, say  $\tilde{f}_n(x)$  where  $\tilde{f}_n$  denotes the approximation of  $n^{th}$  order in a family of approximating functions.

An important aspect of approximating a known or unknown function is the approximation error  $f(x) - \tilde{f}_n(x)$  (sometimes called bias). While uniform upper bounds for the approximation error in the domain  $D$  can be derived for many interpolation methods, these are only a general indication of the precision of the method in cause. In many practical applications, approximation errors can be identified as the source of costs, e.g. in Taguchi's approach to quality optimization. In time series analysis there is a long tradition of evaluating forecast errors according to the costs these incur, where the cost-of-error (loss) function is not necessarily the familiar squared-error loss, see Granger [4].

---

Proceedings of The 4-th International Colloquium "Mathematics in Engineering and Numerical Physics" October 6-8 , 2006, Bucharest, Romania, pp. 68-72.

© Balkan Society of Geometers, Geometry Balkan Press 2007.

Should the function be observed only at some points of its support, or observed with random measurement noise, fitting the desired function to the available, usually experimental, data is the only way to obtain an approximation. Then, following Weiss [7], this function should be fitted under the relevant loss function, i.e. by minimizing mean loss due to approximation error evaluated at each node. Even if  $f$  is completely known, and the parameters of an approximation can be computed analytically (e.g. for a Taylor expansion), these will lead to higher mean loss than fitting.

The main contribution of this note is to study the asymptotic behavior of fitting unknown functions under the relevant loss. We show that, if allowing the approximation order to grow to infinity (this is sometimes called sieve approximation), but at a slower rate than the number of nodes, the fitted parameters converge to their true values, allowing the desired function to be approximated with arbitrary accuracy. If keeping the approximation order constant, the mean cost due to approximation error, although higher than before, is minimized. This adds to a result due to White [8] for the nonlinear least squares fitting procedure.

## 2 The fitting problem

We extend the studied problem to the case where the true value of the function is only given with random additive noise. This allows us to use tools of mathematical statistics that simplify our task, but comes at the cost that all convergence results are stated in probability. However, this is not a serious restriction: the statistical framework reduces to the stated interpolation problem if letting the probability distribution of the noise degenerate to a constant. In this case, convergence in probability turns into standard convergence.

**Assumption 1.** *Let  $y_t = f(x_t) + \varepsilon_t$ , where  $\varepsilon_t$  is independent and identically distributed.*

The *iid* assumption is standard for controlled experiments. With observed data (e.g. in social sciences), this may not be fulfilled and leads to complications, see Christoffersen and Diebold [3].

**Assumption 2.** *Let  $T \rightarrow \infty$  such that the empirical distribution function of the nodes  $x_t$  converges to a proper distribution function.*

This ensures, for instance, that  $x_t$  may be treated as stochastic *iid*, and, more importantly, that  $x_t$  is independent of  $\varepsilon_t$ , see Amemiya [1].

Denote now  $\{\theta_i\}_{i \in \mathbb{N}}$  a sequence of parameters. To ease the exposition, let  ${}_n\theta = (\theta_0, \theta_1, \theta_2, \dots, \theta_n)'$ .

**Assumption 3.** *There is a sequence  $\{\tilde{f}_n\}_{n \in \mathbb{N}}$  of approximating functions,  $\tilde{f}_n : D \times \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ , such that a unique sequence of parameters  $\{\tilde{\theta}_i\}_{i \in \mathbb{N}}$  exists for any function  $f(x) : D \rightarrow \mathbb{R}$  in a given set with  $\tilde{f}_n(x, {}_n\tilde{\theta}) \rightarrow f(x)$  as  $n \rightarrow \infty$  uniformly  $\forall x \in D$ . Further, let  $\tilde{f}_n(x, {}_n\tilde{\theta})$  be Lipschitz continuous w.r.t.  ${}_n\tilde{\theta}$ .*

Taylor polynomials, Hermite polynomials, Tschebyscheff polynomials, Fourier series or Neural Networks can be checked to fulfill Assumption 3.

Granger [4] stated general conditions for loss functions. We require additional regularity conditions, namely continuity and convexity:

**Assumption 4.** *Let  $\mathcal{L}$  be continuous and convex on  $\mathbb{R}$ , increasing on  $\mathbb{R}_+$  and decreasing on  $\mathbb{R}_-$  with  $\mathcal{L}(0) = 0$ .*

Convexity can be dropped at the cost of a more complicated proof, see Remark 1.

**Assumption 5.** *Let  $E(\mathcal{L}(\varepsilon_t - b)) < \infty$ ,  $\forall b \in \mathbb{R}$ , and assume further that  $\arg \min_{b \in \mathbb{R}} E(\mathcal{L}(\varepsilon_t - b)) = 0$ .*

The parameters for the approximation are then estimated from

$${}_n\hat{\theta} = \arg \min_{\theta \in \mathbb{R}^{n+1}} \frac{1}{T} \sum_{t=1}^T \mathcal{L}(y_t - \tilde{f}_n(x_t, \theta)).$$

We address the question of convergence of  ${}_n\hat{\theta}$  to a limiting sequence  $\{\theta_i^*\}_{i \in \mathbb{N}}$  in the following sense

$$\left\| {}_n\hat{\theta} - {}_n\theta^* \right\|_1 \xrightarrow{p} 0 \text{ as } n, T \rightarrow \infty,$$

where  $\|\cdot\|_1$  denotes the  $L_1$  vector norm. If  $\{\theta_i^*\}_{i \in \mathbb{N}} \equiv \{\tilde{\theta}_i\}_{i \in \mathbb{N}}$ , this convergence (stronger than elementwise convergence), together with Assumption 3, implies uniform convergence in probability of  $\tilde{f}_n(x, {}_n\hat{\theta})$  to  $f(x)$ .

### 3 Convergence result

We examine first the case where  $n \rightarrow \infty$ , but slower than  $T$ , i.e.  $n/T \rightarrow 0$ . We show the estimated parameters to converge to the true values,  $\{\tilde{\theta}_i\}_{i \in \mathbb{N}}$ .

**Proposition 1.** *Under Assumptions 1 through 5, it holds as  $n, T \rightarrow \infty$  and  $n/T \rightarrow 0$  that*

$$\left\| {}_n\hat{\theta} - {}_n\tilde{\theta} \right\|_1 \xrightarrow{p} 0.$$

*Proof.* The result is derived in two steps. In the first step, the target function  $T^{-1} \sum_{t=1}^T \mathcal{L}(y_t - \tilde{f}_n(x_t, \theta))$  is shown to converge pointwise in probability as  $n, T \rightarrow \infty$  to a deterministic function minimized only at  $\{\tilde{\theta}_i\}_{i \in \mathbb{N}}$ . In the second, pointwise convergence is shown to imply uniform convergence. Then, using continuity of the arg min operator w.r.t.  $\|\cdot\|_1$  and the supremum norm, the desired result is established. To prove step 1, note that

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathcal{L}(y_t - \tilde{f}_n(x_t, {}_n\theta)) &= \frac{1}{T} \sum_{t=1}^T \mathcal{L}(f(x_t) + \varepsilon_t - \tilde{f}_n(x_t, {}_n\theta)) \\
&= \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\varepsilon_t + f(x_t) - \tilde{f}_n(x_t, {}_n\theta)).
\end{aligned}$$

At  ${}_n\theta = {}_n\tilde{\theta}$ , we have  $f(x_t) - \tilde{f}_n(x_t, {}_n\tilde{\theta}) = o_p(1)$ , with  $o_p(\cdot)$  the correspondent of the Landau symbol in probabilistic terms, and thus

$$\mathcal{L}(\varepsilon_t + f(x_t) - \tilde{f}_n(x_t, {}_n\theta)) = \mathcal{L}(\varepsilon_t) + o_p(1),$$

so

$$\frac{1}{T} \sum_{t=1}^T \mathcal{L}(y_t - \tilde{f}_n(x_t, {}_n\theta)) = \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\varepsilon_t) + o_p(1).$$

A Law of Large Numbers for *iid* variables thus delivers

$$\frac{1}{T} \sum_{t=1}^T \mathcal{L}(y_t - \tilde{f}_n(x_t, {}_n\theta)) \xrightarrow{P} E(\mathcal{L}(\varepsilon_t)).$$

At  ${}_n\theta \neq {}_n\tilde{\theta}$ , it follows that  $\tilde{f}_n(x, {}_n\theta) \rightarrow f^*(x)$ , with  $|f^*(x) - f(x)| > 0 \ \forall x \in D$ . Hence,

$$\begin{aligned}
\frac{1}{T} \sum_{t=1}^T \mathcal{L}(y_t - \tilde{f}_n(x_t, {}_n\theta)) &= \frac{1}{T} \sum_{t=1}^T \mathcal{L}(\varepsilon_t + f(x_t) - f^*(x_t)) + o_p(1) \\
&\xrightarrow{P} E(\mathcal{L}(\varepsilon_t + f(x_t) - f^*(x_t))),
\end{aligned}$$

where the expectation is taken w.r.t. the joint distribution of  $\varepsilon_t$  and  $x_t$ . It was shown by Heike and Demetrescu [5] that

$$E(\mathcal{L}(\varepsilon_t)) < E(\mathcal{L}(\varepsilon_t + u_t)),$$

if  $u_t$  is a random variable, independent of  $\varepsilon_t$ , and non-zero with probability 1, see their Lemma 1. Since, due to Assumption 2,  $f(x_t) - f^*(x_t)$  fulfills these conditions, the limit of the target function is minimized at the true parameters alone. Step 2 is easily completed with the result of Andersen and Gill [2], who showed uniform convergence to follow from pointwise convergence and convexity of the loss function, see their Lemma II.1. The result follows.  $\square$

**Remark 1.** Step 2 of the proof can also be completed by showing stochastic equicontinuity of the sequence of target functions. Uniform convergence then follows due to a result of Newey [6]. This, however, complicates the proof significantly.

Let us now consider the case where  $n$  is fixed. The estimated parameters will not converge to the true values,  $\{\tilde{\theta}_i\}_{i \in \mathbb{N}}$ , but to those parameter values  ${}_n\theta$  that minimize the overall expected loss due to measurement noise and bias  $b(x) = f(x) - \tilde{f}_n(x, {}_n\theta)$ , where the expectation is taken w.r.t. the distribution of  $x_t$  and  $\varepsilon_t$  jointly.

**Proposition 2.** *Under Assumptions 1 through 5, it holds as  $T \rightarrow \infty$  and  $n < \infty$  that*

$${}_n\hat{\theta} \xrightarrow{p} \arg \min_{\theta \in \mathbb{R}^{n+1}} E(\mathcal{L}(\varepsilon_t + b(x_t))).$$

*Proof.* By arguments similar to those in the proof of Proposition 1, the target function converges pointwise to  $E(\mathcal{L}(\varepsilon_t + b(x_t)))$ . The result follows as before.  $\square$

**Remark 2.** *Should  $\varepsilon_t = 0$  with probability 1, the fitted parameters converge to those parameters minimizing the expected bias.*

## 4 Concluding remarks

Asymptotic treatment of the problem of fitting a possibly unknown function under a general cost-of-error (loss) function is provided. Even if analytical approximations are available, fitting delivers better results in terms of overall loss due to approximation error.

## References

- [1] T. Amemiya, *Advanced Econometrics*, Harvard University Press, 1985.
- [2] P.K. Andersen and R.D. Gill, *Cox's regression model for counting processes: a large sample study*, The Annals of Statistics 10 (1982), 1100-1120.
- [3] P.F. Christoffersen and F.X. Diebold, *Optimal prediction under asymmetric loss*, Econometric Theory 13 (1997), 808-817.
- [4] C.W.J. Granger, *Prediction with a generalized cost of error function*, Operational Research Quarterly 20 (1969), 451-468.
- [5] H.-D. Heike and M. Demetrescu (2006), *Forecasting stationary processes under asymmetric loss*, mimeo.
- [6] W.K. Newey, *Uniform convergence in probability and stochastic equicontinuity*, Econometrica 59 (1991), 1161-1167.
- [7] A.A. Weiss, *Estimating time series models using the relevant cost function*, Journal of Applied Econometrics 11 (1996), 539-560.
- [8] H. White, *Consequences and detection of misspecified nonlinear regression models*, Journal of the American Statistical Association 76 (1981), 419-433.

*Authors' addresses:*

Hans-Dieter Heike  
 Statistics and Econometrics, Technical University Darmstadt,  
 Residenzschloss, 64283 Darmstadt, Germany.  
 e-mail: heike@vwl.tu-darmstadt.de

Matei Demetrescu  
 Statistics and Econometric Methods, Goethe-University Frankfurt,  
 Gräffstr. 78 PF 76, 60054 Frankfurt, Germany.  
 e-mail: deme@wiwi.uni-frankfurt.de