

1983

# Some Asymptotic Estimates for Higher-Order Averaging and a Comparison with Iterated Averaging

James Murdock

*Iowa State University*, [jmurdock@iastate.edu](mailto:jmurdock@iastate.edu)

Follow this and additional works at: [http://lib.dr.iastate.edu/math\\_pubs](http://lib.dr.iastate.edu/math_pubs)



Part of the [Mathematics Commons](#)

The complete bibliographic information for this item can be found at [http://lib.dr.iastate.edu/math\\_pubs/114](http://lib.dr.iastate.edu/math_pubs/114). For information on how to cite this item, please visit <http://lib.dr.iastate.edu/howtocite.html>.

---

This Article is brought to you for free and open access by the Mathematics at Iowa State University Digital Repository. It has been accepted for inclusion in Mathematics Publications by an authorized administrator of Iowa State University Digital Repository. For more information, please contact [digirep@iastate.edu](mailto:digirep@iastate.edu).

## SOME ASYMPTOTIC ESTIMATES FOR HIGHER ORDER AVERAGING AND A COMPARISON WITH ITERATED AVERAGING\*

JAMES A. MURDOCK<sup>†</sup>

**Abstract.** Asymptotic estimates for classical higher order averaging are obtained on intervals of length greater than  $O(1/\epsilon)$  when some of the averages vanish. These results are compared with results of Persek using iterated averaging, and the classical methods are found to be more powerful.

**1. Introduction.** We shall be concerned with the  $n$ -dimensional system of differential equations

$$(1) \quad \dot{x} = \epsilon f(x, t, \epsilon) = \epsilon f_1(x, t) + \epsilon^2 f_2(x, t) + \cdots + \epsilon^n f_n(x, t) + \epsilon^{n+1} r_{n+1}(x, t, \epsilon)$$

where  $f$  is  $2\pi$ -periodic in  $t$  and  $\epsilon$  is a small positive real number. For such systems the traditional  $n$ th order averaging method, as described for instance in Perko [2], yields approximate solutions which retain accuracy  $O(\epsilon^n)$  on time intervals of length  $O(1/\epsilon)$ . It is of interest whether one can "trade off" some of this accuracy for validity on a longer time interval; that is to say, we may ask whether the same approximate solution retains the decreased accuracy  $O(\epsilon^{n-j})$  on the expanded time interval of length  $O(1/\epsilon^{1+j})$  for certain integers  $j$ . Our first theorem (§2) asserts that this is true for  $j=0, 1, \dots, l-1$  if the  $n$ th order averaged system corresponding to (1) begins with the term of order  $\epsilon^l$ , in other words takes the form

$$(2) \quad \dot{z} = \epsilon^l g_l(z) + \cdots + \epsilon^n g_n(z).$$

The proof involves no new methods.

Persek [3] has defined a method which he calls "iterated averaging" which under certain conditions approximates system (1) by a system of the form

$$(3) \quad \dot{z} = \epsilon^l h_l(z).$$

He then proves that solutions of (3) approximate those of (1) to order  $O(\epsilon)$  on an interval of length  $O(1/\epsilon^l)$ . It is evident that (3) has the same form as (2) in the case  $n=l$ , although it is not clear a priori whether the function  $h_l$  constructed by Persek coincides with the classical  $g_l$  in this case. Nevertheless Persek's estimate of the accuracy of (3) coincides with our estimate for (2) if  $n=l$  and if  $j$  is taken to be  $l-1$ . This prompts a comparison between  $g_l$  and  $h_l$ , which we carry out in §3 for the case  $l=2$ . Briefly the result is that  $g_2$  and  $h_2$  coincide when both are defined, but that the defineability of  $h_2$  depends upon a hypothesis which is unnecessary under the traditional method. Thus Persek's result is (at least for the case  $l=2$ ) a special case of ours (in §2). For  $l>2$  it would be tedious to compare  $h_l$  and  $g_l$ , but it is again apparent that  $g_l$  is always defined (and thus equation (2) exists provided that  $g_1$  through  $g_{l-1}$  vanish), whereas  $h_l$  is defined only if  $h_1$  through  $h_{l-1}$  vanish and in addition a further condition is satisfied, which does not correspond to any condition in the traditional theory.

\*Received by the editors November 13, 1981.

<sup>†</sup>Department of Mathematics, Iowa State University, Ames, Iowa 50011.

**2. Asymptotic estimates.** It is shown in the classical theory of averaging that there exists a coordinate transformation

$$(4) \quad x = u(y, t, \epsilon) = y + \epsilon u_1(y, t) + \dots + \epsilon^n u_n(y, t)$$

$2\pi$ -periodic in  $t$  and carrying (1) to

$$(5) \quad \dot{y} = \epsilon g_1(y) + \dots + \epsilon^n g_n(y) + \epsilon^{n+1} R_{n+1}(y, t, \epsilon).$$

The transformation (4) is not unique and is usually normalized either by requiring that the average of each  $u_k$  vanish, or by requiring that each  $u_k$  vanish for  $t=0$ . The latter is called the stroboscopic method and has the advantage that (4) reduces to  $x=y$  at  $t=0$  and at all stroboscopic times  $t=2\pi, 4\pi, \dots$ . Associated with (5) is the truncated system

$$(6) \quad \dot{z} = \epsilon g_1(z) + \dots + \epsilon^n g_n(z)$$

which is wholly autonomous. Let  $x(t, \epsilon), y(t, \epsilon), z(t, \epsilon)$  denote solutions of (1), (5), and (6) defined in an interval  $0 \leq \epsilon \leq \epsilon_0$  whose initial conditions are related by  $x(0, \epsilon) = u(y(0, \epsilon), 0, \epsilon)$  and  $y(0, \epsilon) = z(0, \epsilon)$ ; note that in the stroboscopic case this reduces to  $x(0, \epsilon) = y(0, \epsilon) = z(0, \epsilon)$ . The classical method of averaging proposes to construct from  $z(t)$  an approximation to  $x(t)$ . Since  $z(t)$  ought to be close to  $y(t)$ , and since  $y(t)$  is related to  $x(t)$  by (4), in view of the fact that (4) is assumed to hold at  $t=0$ ,  $x(t)$  should be well approximated by  $u(z(t, \epsilon), t, \epsilon)$ , an expression which is called the *improved  $n$ th approximation* to  $x(t)$ . It turns out, however, that there is no loss in asymptotic accuracy if the term of order  $n$  in (4) is omitted in forming the approximation (it must not be omitted, of course, in transforming (1) into (5)). Therefore the  *$n$ th approximation* to  $x(t)$  is defined by

$$(7) \quad X(t, \epsilon) = \hat{u}(z(t, \epsilon), t, \epsilon),$$

where

$$\hat{u}(z, t, \epsilon) = z + \epsilon u_1(z, t) + \dots + \epsilon^{n-1} u_{n-1}(z, t).$$

**THEOREM.** *Under the above hypotheses there exist positive constants  $c$  and  $c_0$  such that  $|x(t, \epsilon) - X(t, \epsilon)| \leq c_0 \epsilon^n$  for  $0 \leq t \leq c/\epsilon$  and  $0 \leq \epsilon \leq \epsilon_0$ . Under the additional hypothesis that  $g_1$  through  $g_{l-1}$  vanish, so that (6) takes the form (2), there exist for each  $j=0, \dots, l-1$  positive constants  $c_j$  such that  $|x(t, \epsilon) - X(t, \epsilon)| \leq c_j \epsilon^{n-j}$  for  $0 \leq t \leq c/\epsilon^{1+j}$ .*

*Proof.* We prove the second assertion. The first, which is classical, is included by taking  $l=1, j=0$ .

Let  $K$  be a closed ball centered at  $x(0, 0)$  of radius  $R$  sufficiently large that  $x(0, \epsilon)$  and  $y(0, \epsilon) = z(0, \epsilon)$  are contained in the concentric ball of radius  $R/2$  for  $0 \leq \epsilon \leq \epsilon_0$ . Since  $K$  is compact there exists in view of (1), (2) and (5) a constant  $M$  such that  $d|y|/dt$  and  $d|z|/dt$  are less than  $M\epsilon^l$  for  $0 \leq \epsilon \leq \epsilon_0$  as long as  $y$  and  $z$  remain in  $K$ . Since  $|y|$  and  $|z|$  must drift by at least  $R/2$  from their initial positions in order to leave  $K$ , and since in time  $t$  they can drift at most  $M\epsilon^l t$  while in  $K$ , it follows that there exists a constant  $c > 0$  such that  $y$  and  $z$  remain in  $K$  for  $0 \leq t \leq c/\epsilon^l$ . Hence on this interval (the longest considered in the theorem) one may use, for all functions of  $y$  and  $z$ , their upper bounds and Lipschitz constants on  $K$ .

Letting  $\rho = |y(t, \epsilon) - z(t, \epsilon)|$  one immediately finds from (2) and (5) Lipschitz constants  $L_1, \dots, L_n$  and a bound  $B$  such that  $d\rho/dt \leq (\epsilon^l L_1 + \dots + \epsilon^n L_n)\rho + \epsilon^{n+1} B$ . Solving this linear differential inequality with initial condition  $\rho=0$  yields  $\rho \leq \epsilon^{n+1} B \delta^{-1} (e^{\delta t} - 1)$  where  $\delta = \epsilon^l L_1 + \dots + \epsilon^n L_n$ . Estimating the factor  $e^{\delta t} - 1$  requires some care. The tempting answer  $e^{\delta t} - 1 = O(\delta t) = O(\epsilon^l t)$  is correct for the time intervals with which we

are concerned, but this requires proof. Namely  $e^x - 1 = x(1 + x/2! + \dots) = x\phi(x)$  and hence for  $x$  in any bounded interval there is a constant  $k$  such that  $e^x - 1 \leq kx$ ; if  $x = \delta t$  and  $0 \leq t \leq c/\epsilon^l$  and  $0 < \epsilon \leq \epsilon_0$  then  $x$  is bounded and  $e^{\delta t} - 1 \leq k\delta t$ . (On longer intervals  $e^{\delta t} - 1$  can approach infinity faster than  $\delta t$ .) Since  $\delta^{-1} = O(1/\epsilon^l)$  we have  $\rho = O(\epsilon^{n+1}t)$  as long as  $y$  and  $z$  remain in  $K$ , that is, at least on the interval  $0 \leq t \leq c/\epsilon^l$ . It follows that on any interval  $0 \leq t \leq c/\epsilon^{1+j}$ ,  $j=0, \dots, l-1$ , one has  $|y(t, \epsilon) - z(t, \epsilon)| = \rho = O(\epsilon^{n-j})$ ; that is, this quantity is bounded by a constant times  $\epsilon^{n-j}$  for  $0 \leq t \leq c/\epsilon^{1+j}$  and  $0 \leq \epsilon \leq \epsilon_0$ . Now using the Lipschitz constant for  $u$  on  $K$  one finds  $|x(t, \epsilon) - u(z(t, \epsilon), t, \epsilon)| = |u(y(t, \epsilon), t, \epsilon) - u(z(t, \epsilon), t, \epsilon)| = O(\epsilon^{n-j})$  on the same interval. But

$$|u(z(t, \epsilon), t, \epsilon) - X(t, \epsilon)| = |u(z(t, \epsilon), t, \epsilon) - \hat{u}(z(t, \epsilon), t, \epsilon)| = O(\epsilon^n)$$

for all time; adding the last two estimates proves the theorem. Q.E.D.

In the proof it is seen that the final term in the transformation  $u$  is unnecessary in constructing  $X(t)$  because the error committed by leaving it out is of the same order as the error already present. By the same reasoning we see that for  $j > 0$ , where the possible accuracy is at most  $O(\epsilon^{n-j})$ , we may omit  $j$  additional terms from  $u$  in forming  $X$ . In particular, in the case  $n = l$ ,  $j = l - 1$ , it is not necessary to use  $u$  at all and we obtain

**COROLLARY.** *When (2) reduces to  $\dot{z} = \epsilon^l g_l(z)$ , there exist positive constants  $c$  and  $c_0$  such that  $|x(t, \epsilon) - z(t, \epsilon)| < c_0 \epsilon$  for  $0 \leq t \leq c/\epsilon^l$ .*

This form of the theorem is the one most directly comparable to the work of Persek. The comparison is carried out in the next section.

**3. Comparison of two averaging methods.** In order to calculate  $g_2(z)$  it is necessary to recall how (4) and (5) are constructed. It is clear a priori that any transformation of the form (4) carries (1) into a system of the form

$$(8) \quad \dot{y} = \epsilon g_1(y, t) + \dots + \epsilon^n g_n(y, t) + \epsilon^{n+1} R_{n+1}(y, t, \epsilon).$$

From (1), (4), and (8) one calculates that the  $f$ 's,  $u$ 's, and  $g$ 's are related by

$$(9) \quad \begin{aligned} \frac{\partial u_1}{\partial t}(y, t) &= f_1(y, t) - g_1(y, t), \\ \frac{\partial u_2}{\partial t}(y, t) &= \left\{ f_2 + \frac{\partial f_1}{\partial y} u_1 - \frac{\partial u_1}{\partial y} g_1 \right\}_{(y, t)} + g_2(y, t) \end{aligned}$$

with similar equations for the higher  $u_n$ 's. (Briefly, to obtain (9) differentiate (4), insert (5) and compare this with the result of inserting (4) into (1).) It is clear that (9) admits solutions for  $u_1$  and  $u_2$  which are periodic in  $t$  if and only if the right-hand sides have zero mean value. Now to achieve (5) the  $g_i$  must be independent of  $t$ ; thus (9) dictates that  $g_1(y)$  must be the average of  $f_1(y, t)$  and  $g_2(y)$  must be the average of the expression in braces.

We wish to consider the case in which the averaged system takes the form (2) with  $l=2$ . Thus we now assume  $g_1(y) = 0$ , which is to say that the average of  $f_1(y, t)$  vanishes. In this case  $u_1(y, t) = \int_a^t f_1(y, s) ds$  with  $a$  arbitrary ( $a=0$  gives the stroboscopic method). Inserting this into the expression in braces and averaging gives the following formula for  $g_2$ , in which the dependence upon the choice of  $a$  is made explicit:

$$(10) \quad g_2(y, a) = \frac{1}{2\pi} \int_0^{2\pi} \left\{ f_2(y, t) + \frac{\partial f_1}{\partial y}(y, t) \int_a^t f_1(y, s) ds \right\} dt.$$

Downloaded 07/12/17 to 129.186.176.188. Redistribution subject to SIAM license or copyright; see http://www.siam.org/journals/ojsa.php

Persek's function  $h_2(z)$  is defined as follows in our notation (compare his  $\bar{E}^{(2)}$ , [3, p. 416]). First assume the average of  $f_1(y, t)$  vanishes, as we have done above. Next define

$$(11) \quad H_2(z, \tau) = \frac{1}{2\pi} \int_{\tau}^{\tau+2\pi} \left\{ f_2(z, t) + \frac{\partial f_1}{\partial y}(y, t) \int_{\tau}^t f_1(y, s) ds \right\} dt.$$

If the latter expression is independent of  $\tau$ , it is defined to be  $h_2(z)$ ; if it is not independent of  $\tau$ ,  $h_2(z)$  is not defined.

Now it is clear that the bracketed expression in (10) is periodic in  $t$ , since we have assumed that  $f_1$  has zero mean value. Therefore  $\int_0^{2\pi}$  in (10) may be replaced by  $\int_{\tau}^{\tau+2\pi}$  for any  $\tau$ . The only remaining difference between (10) and (11) is that  $a$  in (10) is replaced by  $\tau$  in (11). Thus we see that

$$(12) \quad H_2(z, \tau) = g_2(z, \tau).$$

Thus the sole difference between Persek's average (for  $n=l=2$ ) and ours is that Persek must assume (11) independent of  $\tau$ , whereas the corresponding quantity  $a$  in (10) enters as an arbitrary constant and requires no additional assumptions. It is clear that the assumption that (11) is independent of  $\tau$  is a very strong assumption and one which gains no advantage.

With regard to higher order terms the following situation obtains. Our  $g_k$  is always definable, and is not unique but rather depends upon the choices of integration constants in solving for  $u$ . On the other hand  $h_k$  is only defined if two conditions are met:  $h_1, \dots, h_{k-1}$  must be defined and vanish, and a certain function  $H_k(z, \tau)$  (which Persek calls  $\bar{E}^{(k)}$ ) must be independent of  $\tau$ . The presence of the latter restriction indicates that there are likely to be many cases in which (2) takes the form  $\dot{z} = \varepsilon^l g_l(z)$  and yet (3) cannot be formulated. It seems reasonable to conjecture (based on the case  $l=2$ ) that  $h_l$  exists precisely when  $g_l$  is unique (i.e., independent of the choices made in  $u$ ) and that in this case  $h_l = g_l$ . Proof of such a theorem, if true, would involve notational difficulties but might be attempted (if it were considered important) by the use of the operators constructed by Musen [1] for use in averaging methods, based upon a formula of St. Faa de Bruno.

#### REFERENCES

- [1] P. MUSEN, *On the high order effects in the methods of Krylov-Bogoliubov and Poincaré*, J. Astronaut. Sci. 12 (1965), pp. 129-134.
- [2] L. M. PERKO, *Higher order averaging and related methods for perturbed periodic and quasi-periodic systems*, SIAM J. Appl. Math., 17 (1968), pp. 698-723.
- [3] S. C. PERSEK, *Hierarchies of iterated averages for systems of ordinary differential equations with a small parameter*, SIAM J. Math. Anal., 12 (1981), pp. 413-420.