

An algorithm and computer program for the calculation of envelope curves

Marjorie McClain, Albert Feldman, David Kahaner, and Xuantong Ying^{a)}
National Institute of Standards and Technology, Gaithersburg, Maryland 20899

(Received 3 January 1990; accepted 30 April 1990)

A procedure has been developed to calculate numerically the envelope functions of an oscillatory curve. The method has been shown to be applicable to optical transmission data, but it is general enough to be used for many other data sets. The program is available on request.

INTRODUCTION

In many instances, experimental data consist of a series of oscillations bounded by upper and lower curves which are envelopes. Such behavior is frequently observed when interference effects are superimposed on a slowly varying trend. By knowing the envelope functions one may be able to deduce physical properties giving rise to the observed phenomenon. An example of this is the envelope method used by Manifacier *et al.*¹ to obtain the optical constants and thicknesses of weakly absorbing films from transmission spectra. It would therefore be useful if a simple automated procedure could be used to obtain the envelope functions, especially in numerical form. In this paper we describe an iterative method for calculating the envelope curves of a given set of oscillatory data. In an example, we analyze a set of transmission data on the basis of the paper by Manifacier *et al.*¹ and compare the results with a manual method.

I. THE PROCEDURE

The envelopes of a given oscillatory function $T(x)$ are two smooth curves that in some sense represent the maximum and minimum values of $T(x)$. By "smooth" it is meant here that the envelopes have very few inflection points compared with $T(x)$; often they have no inflection points at all. The envelopes are constrained to lie tangent to $T(x)$ and not to cross it. The top envelope $T_{\max}(x)$ lies above the function, i.e., $T_{\max}(x) \geq T(x)$, while the bottom envelope $T_{\min}(x)$ lies below the function, i.e., $T_{\min}(x) \leq T(x)$. Figure 1 shows a schematic representation of an oscillatory function $T(x)$ and its envelopes $T_{\max}(x)$ and $T_{\min}(x)$. Note that the points of tangency are not the same as the critical points (i.e., the local extrema) of $T(x)$. The tangent points often lie near the critical points, but this is not always the case; see, for example, the left-most tangent point in Fig. 1.

The following steps outline a procedure for calculating envelope curves:

- (1) Smooth the given data.
- (2) Estimate the locations of the upper and lower tangent points.
- (3) Interpolate a curve through the estimated upper tangent points and another through the estimated lower tangent points.
- (4) If no points on the upper curve lie below the smoothed data and no points on the lower curve lie above the smoothed data, then stop. The upper curve is the top envelope T_{\max} , and the lower curve is the bottom envelope T_{\min} .
- (5) Otherwise, improve the estimates of the tangent point locations and return to step 3.

In what follows, the calculation of the top envelope is discussed in detail. The calculation of the bottom envelope is exactly analogous.

The only input required by the envelope algorithm is the set of data points (x,y) whose envelope is to be computed. Due to the likely presence of small errors in the data, the original set of data points is not used directly but is replaced

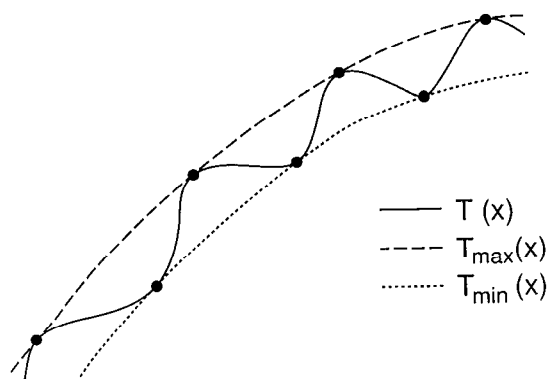


FIG. 1. Schematic representation of an oscillatory function and its envelopes. The dots represent the points of tangency of the function with its envelopes.

^{a)} Guest scientist from Fudan University, Department of Physics, Laboratory of Laser Physics and Optics, Shanghai, People's Republic of China.

by a smoother set (step 1 above). The smoothed data points are obtained by evaluating a spline function $T(x)$ that is computed to be a least-squares fit to the original data. $T(x)$ is required to lie within a specified error tolerance ϵ of each original data point; that is, for each data point (x, y) , $|T(x) - y| \leq \epsilon$. Two separate smoothing functions are employed in the course of the envelope calculation: a rough initial smoothing $T_0(x)$ that is used in the first estimation of the tangent point locations and a more accurate final smoothing $T(x)$ that is used in determining the envelope curve.

The envelope algorithm starts by estimating the locations of the points where $T_{\max}(x)$ is tangent to $T_0(x)$ (step 2 above). The obvious choices for these tangent point estimates are the local maxima of $T_0(x)$. However, as noted above, not every tangent point is near a maximum point, so a more generally applicable estimation method is employed. The algorithm first finds the intervals where $T_0(x)$ opens downward, i.e., where $T_0''(x) < 0$. It then places a tangent point estimate halfway between the end points of each of these intervals, as shown in Fig. 2. The end points are the inflection points of $T_0(x)$, i.e., points where $T_0''(x) = 0$. (The second derivative information needed for determining the inflection points and the direction of curvature is obtained during the smoothing process in step 1.)

Downward-opening intervals containing the first or last data point ("end-point intervals") require special treatment. In some cases, the data for these intervals might be incomplete; that is, if the data could be extended, $T_0(x)$ would continue to open downward for some interval beyond the original end point. The true tangent point for such an interval might then lie outside the range of the original data, as shown in Fig. 3. The envelope algorithm does not attempt to estimate such tangent points. It will produce a tangent point estimate for an end-point interval only if it contains a point where $T_0'(x) = 0$. The tangent point estimate will then be set to be this local maximum point, since presumably the true tangent point will lie nearby, within the range of the data. Recall again that not every tangent point is near a maximum point, so this method for handling end-point intervals will occasionally miss valid tangent points, but it will not include invalid ones.

The initial estimation of the tangent points is the most

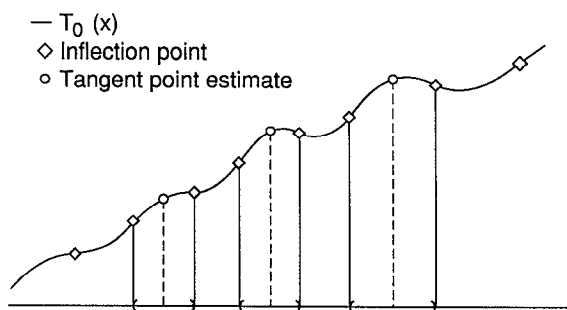


FIG. 2. A method for obtaining initial estimates of tangent points. A tangent estimate is set at the midpoint of each downward-opening interval bounded by two inflection points.

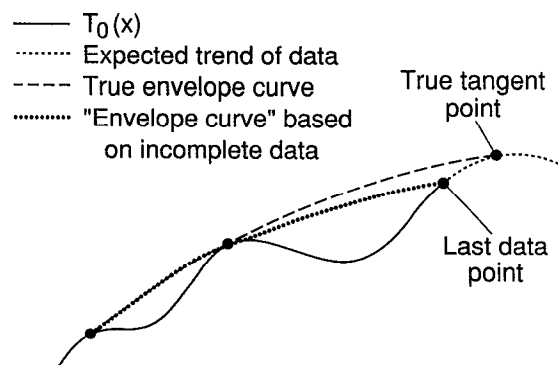


FIG. 3. A problem with estimating tangent points near the ends of the data. In this figure, the right-most downward-opening interval of the data curve is incomplete. If it could be extended (as the dotted curve indicates), the true tangent point would be seen to lie beyond the end of the original data. Using any point in the existing data as a tangent point estimate would produce an "envelope curve" very different from the true envelope curve.

critical part of the envelope calculation, and it depends heavily on the initial smoothing of the data. If the smoothed curve $T_0(x)$ does not match the data closely enough, there will be too few inflection points and hence too few tangent points. On the other hand, if $T_0(x)$ matches the data too closely, there are likely to be many undesirable inflection points, and these will produce spurious tangent point estimates, as shown in Fig. 4. It may be necessary to experiment with various error tolerances ϵ before finding a $T_0(x)$ that leads to appropriate tangent point estimates.

Next, a first approximation to the envelope curve is obtained by interpolating a smooth curve through the estimated tangent points (step 3 above). The interpolated curve $I(x)$ is required to preserve monotonic behavior; that is, if the tangent points increase (or decrease) with x , the resulting envelope curve does also. If the tangent points are not monotonic, the interpolation method forces an extreme point in the envelope curve at any point where the direction of monotonicity changes.

The first approximation to the envelope curve will generally not be truly tangent to the smoothed data curve, unless the initial tangent point estimates happen to be exactly correct. In most cases, the interpolated curve will cut

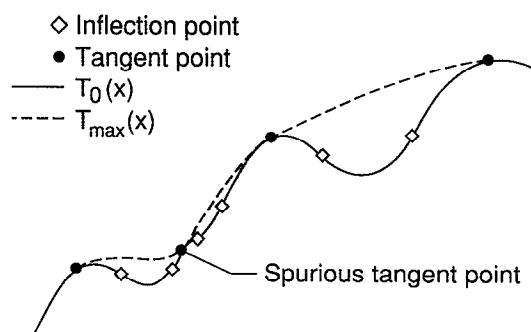


FIG. 4. A spurious tangent point, caused by a small wiggle in the smoothed data curve. This leads to an unacceptable envelope curve.

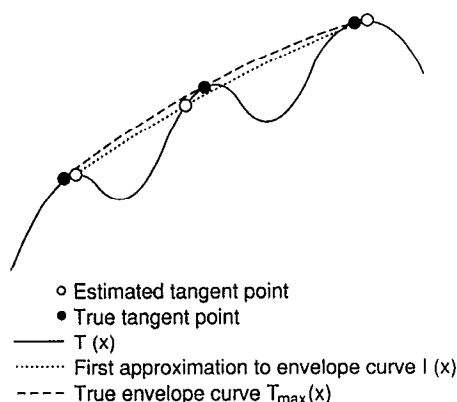


FIG. 5. The first approximation to the envelope curve usually lies below the true envelope.

below the peaks of the data curve, as shown in Fig. 5. Around each of the estimated tangent points, there will be an interval where $I(x) < T(x)$. (The estimated tangent point will actually be one of the end points of the interval.) The points of tangency of the true envelope $T_{\max}(x)$ with $T(x)$ must lie somewhere on these intervals. The algorithm proceeds (step 5 above) by taking the midpoints of each of these intervals as new estimates for the tangent points, as shown in Fig. 6.

When checking for envelope points that lie below $T(x)$, it is generally desirable to use a smoothing fit that tightly matches the original data. This is especially necessary in the regions closely surrounding the tangent points, in order that the final envelope curve will be as accurate as possible. The initial smoothing function $T_0(x)$ used to obtain the first tangent point estimates may not have sufficient accuracy to serve as the final smoothing, so it is generally necessary to try another spline fit $T(x)$ with a smaller error tolerance ϵ in order to produce the desired accuracy in the final envelope curve.

When the new tangent point estimates have been computed, a new curve $I(x)$ is interpolated through them, which should lie closer to the true envelope. This process of

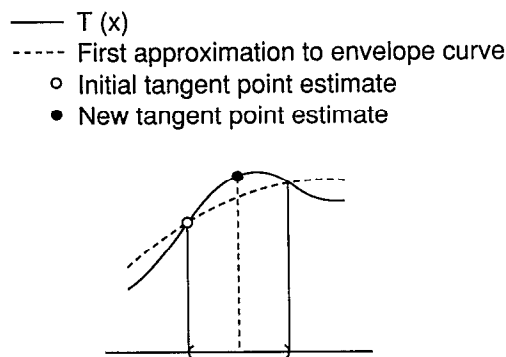


FIG. 6. Obtaining an improved tangent point estimate. The first approximation to the envelope curve cuts off the peaks of the data curve. A new tangent point estimate is set at the midpoint of each cut-off interval.

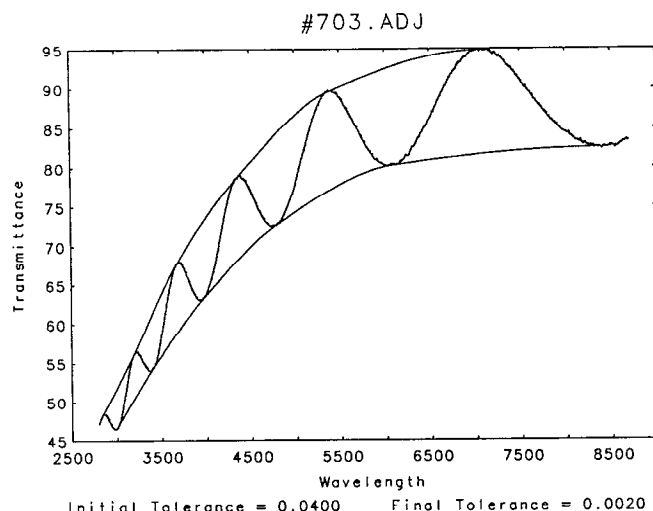


FIG. 7. Transmittance versus wavelength of a mixed yttria-silica film after subtraction of reflection from the air-silica surface of the substrate. This figure, showing experimental data, smoothed data, and fitted envelope curves, was produced directly by the envelope program.

finding new estimates for the tangent points and interpolating a curve through them is repeated until either no points on the interpolated curve lie beneath the data curve (step 4 above) or until a maximum of 20 iterations have been performed. The resulting interpolated curve $I(x)$ should lie very close to the true envelope $T_{\max}(x)$.

We have implemented this envelope algorithm as a FORTRAN 77 program running on an IBM PC or compatible computer. To smooth the data, we use the least-squares spline-fitting subroutine EFC,² augmented by a subroutine that automatically generates break points for the spline (i.e., the points where the polynomial pieces of the spline are joined). To interpolate a monotonic function through the tangent estimates, we use the subroutine PCHIM.³ Both of these subroutines are part of SLATEC,⁴ a public domain mathematical software library available from Argonne National Laboratory. The main program allows the user to request a data set and to enter the tolerance factors to be used for the data smoothing. The data and the resulting envelope curves are then plotted on the screen. The user has several options: (1) redo the envelope calculation with different error tolerances; (2) produce a high-quality hard copy of the plot; or (3) save the envelope data and tangent point locations in a file.

A typical wall-clock running time for the program on an IBM PC/AT is about 1.75 min for a data set containing 800 points. (This time does not include plotting and saving the results.) Approximately 1.25 min are spent on smoothing the data, while 0.50 min are spent on calculating the envelopes.

II. AN APPLICATION

Figure 7 shows the optical transmission of film composed of a mixture of yttria and silica deposited on a fused silica substrate. Shown superimposed on the data are the enve-

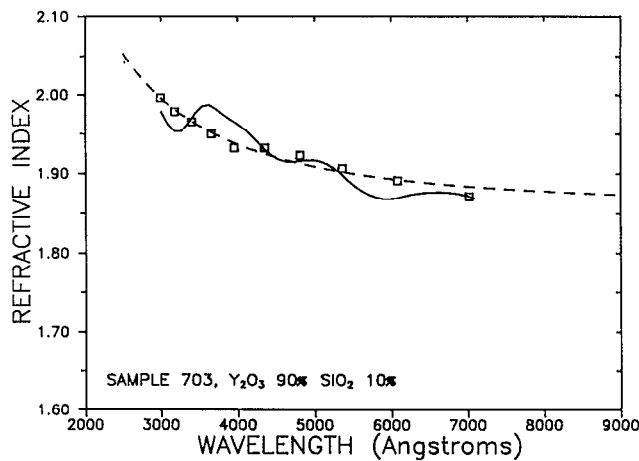


FIG. 8. Refractive index versus wavelength of a mixed yttria-silica film. The squares represent results obtained by a prior manual analysis, the solid curve represents results calculated from the fitted envelope curves, and the dashed curve represents a fit to the manual data.

lope functions computed by the procedure described above. The method of Manifacier *et al.*¹ has been used to obtain the refractive index n , the absorption coefficient α , and the thickness t , of the film from the envelope functions.

The transmittance T of a weakly absorbing film on a transparent substrate can be represented by

$$T = \frac{n_0 n^2 n_s A}{C_1^2 + C_2^2 A^2 + 2C_1 C_2 A \cos(4\pi n t / \lambda)}, \quad (1)$$

where $C_1 = (n + n_0)(n_s + n)$, $C_2 = (n - n_0)(n_s - n)$, $A = \exp(-4\pi k t / \lambda) = \exp(-\alpha t)$, n_0 is the refractive index of the ambient (air), n_s is the refractive index of the substrate, k is the imaginary part of the refractive index of the film, and λ is the wavelength of the radiation. In general, each of the parameters except t are functions of λ . Equation (1) is an oscillatory function with envelope curves given by

$$T_{\max} = n_0 n^2 n_s A / (C_1 + C_2 A)^2 \quad (2)$$

and

$$T_{\min} = n_0 n^2 n_s A / (C_1 - C_2 A)^2, \quad (3)$$

as can be seen by taking the cosine to be $+1$ or -1 . While there is no guarantee that the envelope algorithm described in this paper will produce the curves given by Eqs. (2) and (3), it is believed to give a close approximation. The procedures for calculating n , α (or k), and t from the envelope functions and the points of tangency will not be repeated here, as they are discussed in Manifacier *et al.* The results for the refractive index data are shown in Fig. 8, where n is plotted as a function of λ . The figure also shows the values obtained by an earlier manual procedure,⁵ and it can be seen that the two methods agree reasonably well. (Note that the vertical scale is expanded and does not begin at zero.) A similar plot for the absorption coefficient data, not shown here, exhibits even more precise agreement between the two methods. The value obtained for thickness was $0.58 \mu\text{m}$, which agrees reasonably well with a mechanical thickness measurement of $0.60 \mu\text{m}$ and a value from the manual data analysis procedure of $0.56 \mu\text{m}$.

III. HOW TO OBTAIN COPIES OF THE PROGRAM

Copies of the envelope program may be obtained at no charge by contacting the first author (email: marje@fsl.cam.nist.gov). The program will be provided as an executable file of size approximately 307 kbytes, along with a sample data file, on a single $3\frac{1}{2}$ -in. or $5\frac{1}{4}$ -in. DOS diskette. Source code (about 4000 lines of FORTRAN), as well as a graphics library file referenced by the main program, will be included on additional disks if specifically requested.

REFERENCES

1. J. C. Manifacier, J. Gasiot, and J. P. Fillard, "A simple method for the determination of the optical constants n , k and the thickness of a weakly absorbing film," *J. Phys. E.: Sci. Instrum.* **9**, 1002 (1976).
2. R. J. Hanson, *Constrained Least Squares Curve Fitting to Discrete Data Using B-Splines, a User's Guide*, Sandia Laboratories Technical Report SAND-78-1291 (1978).
3. F. N. Fritsch and R. E. Carlson, "Monotone piecewise cubic interpolation," *SIAM J. Num. Anal.* **17**, 238 (1980).
4. B. L. Buzbee, "The SLATEC common mathematical library," in *Sources and Development of Mathematical Software*, edited by W. R. Cowell (Prentice-Hall, Englewood Cliffs, NJ, 1984).
5. A. Feldman, X. T. Ying, and E. N. Farabaugh, "Optical properties of mixed yttria-silica films," *Appl. Opt.* **28**, 5229 (1989).