diation from an isolated system, in terms of the behavior of the fields near null infinity. The Bondi mass was of importance in proving that gravity waves carry away mass-energy: the Bondi mass decreases monotonically whenever gravitational radiation is emitted. A modification of Witten's proof (see Horowitz and Strominger[10]) has demonstrated the positivity of the Bondi mass.

Bondi et al.[9] give an explicit expression for a quantity called the mass aspect of a mass $M$ (a Schwarzschild metric) moving with constant velocity. The mass aspect is a function of the observer's polar angle, measured from the direction of motion of the moving object. The average of the mass aspect over all solid angle is computed and equals $\gamma M$.

However, there is no direct conflict between the value of the Bondi mass ($\gamma M$) for a moving Schwarzschild solution and our result, Eq. (10). The Bondi mass equals $\gamma M$ only when averaged over all solid angle. Our problem, a Schwarzschild object passing by a test particle, is quite different: most of the scattering may take place when the test particle is closest to the Schwarzschild object and the value of the polar angle is near 90°.

[a] Current address: Department of Physics, Southwest Texas State University, San Marcos, TX 78666.

[b] Current address: Department of Physics, Texas A & M University, College Station, TX 77843.

[1] H. Goldstein, Classical Mechanics (Addison-Wesley, London, 1980), 2nd ed., pp. 94–99.

[2] R. A. Matzner and Y. Nutku, Proc. R. Soc. A336, 285 (1974).

[3] S. Weinberg, Gravitation and Cosmology (Wiley, New York, 1972), p. 189.

[4] R. Arnowitt, S. Deser, and C. W. Misner, Phys. Rev. 118, 1100 (1960).

[5] R. Schoen and S.-T. Yau, Commun. Math. Phys. 65, 45 (1979).

[6] R. Schoen and S.-T. Yau, Phys. Rev. Lett. 43, 1457 (1979).

[7] R. Schoen and S.-T. Yau, Commun. Math. Phys. 79, 231 (1981).

[8] E. Witten, Commun. Math. Phys. 80, 381 (1981).

[9] H. Bondi, M. G. van der Burg, and A. W. K. Metzner, Proc. R. Soc. A269, 21 (1962).

[10] G. T. Horowitz and A. Strominger, Phys. Rev. D 27, 2793 (1983).

# Simple examples of correlations in error propagation

John R. Taylor

Department of Physics, University of Colorado, Boulder, Colorado 80309

This paper offers two simple examples of pairs of variables which are correlated, and whose covariance is easily evaluated. Examples are given of error propagation in which these correlations play an important role.

## I. INTRODUCTION

One of the most important aspects of the error analysis that we teach in introductory physics laboratories is error propagation. For example, we teach our students that if they measure two independent variables $s$ and $t$, with uncertainties $\sigma_s$ and $\sigma_t$, and use their measurements to calculate some quantity $q(s,t)$ that depends on $s$ and $t$, then the uncertainty $\sigma_q$ in their answer is given by the quadratic sum

$$\sigma_q^2 = \left(\frac{\partial q}{\partial s}\right)^2 \sigma_s^2 + \left(\frac{\partial q}{\partial t}\right)^2 \sigma_t^2. \tag{1}$$

Unfortunately, as most texts emphasize, this formula is really justified only if the measurements of $s$ and $t$ are uncorrelated. This leaves the students (and teachers) in some confusion as to when they should use Eq. (1).

In more advanced courses and books it is shown that, whether or not $s$ and $t$ are correlated, $\sigma_q$ is given by the formula

$$\sigma_q^2 = \left(\frac{\partial q}{\partial s}\right)^2 \sigma_s^2 + 2\frac{\partial q}{\partial s}\frac{\partial q}{\partial t}\sigma_{st} + \left(\frac{\partial q}{\partial t}\right)^2 \sigma_t^2, \tag{2}$$

where $\sigma_{st}$ is the covariance of $s$ and $t$. If $s$ and $t$ are uncorrelated, then $\sigma_{st} = 0$ and (2) gives back (1). If $\sigma_{st} \neq 0$ then, depending on the signs of $\sigma_{st}$ and of the two derivatives, the correct value of $\sigma_q$ as given by (2) can be larger or smaller than that given by (1).

Unfortunately it is rare in an elementary or intermediate laboratory that one has the opportunity to use the formula (2), because one seldom knows the covariance $\sigma_{st}$. The most that is generally done with (2) is to derive an upper bound on $\sigma_q$, as follows: According to the Schwartz inequality,[1] $|\sigma_{st}| \leqslant \sigma_s \sigma_t$. Therefore, it follows from (2) that

$$\sigma_q \leqslant \left|\frac{\partial q}{\partial s}\right|\sigma_s + \left|\frac{\partial q}{\partial t}\right|\sigma_t. \tag{3}$$

That is, the uncertainty $\sigma_q$ is never worse than the linear sum (3).

To illustrate the relative merits of Eqs. (1) to (3), I describe here two elementary examples where the covariance that appears in (2) is clearly nonzero and can be explicitly evaluated. Further, when one uses the correct formula (2) one can get answers that are strikingly different from both (1) and (3).

In Sec. II, I describe the strong negative correlation of two angles measured on certain bubble-chamber photographs, and in Sec. III the correlation of the least-squares coefficients $A$ and $B$. I do not claim that either example is previously unknown; quite the contrary, the correlation of the least-squares coefficients is certainly well known to people who use least-squares fitting regularly, and is described in several advanced textbooks.[2] However, I have not found any simple examples of correlated variables (in the context of error propagation) in an elementary text.[3] From several inquiries I have received in connection with my own "Error Analysis" (Ref. 1), it appears to me that several teachers and students would find the following examples helpful.

## II. TWO CORRELATED ANGLES

My first example of two correlated variables comes from the introductory modern physics laboratory, Physics 215, which is taken by sophomores at the University of Colorado. In this laboratory, students are given bubble-chamber photographs of the production and decay of neutral $K$ mesons. The events are chosen to lie in the plane of the photograph to avoid problems with three-dimensional geometry. A typical event is sketched in Fig. 1, which shows the process

$$K^+ + n \rightarrow K^0 + p \qquad (4)$$

followed by

$$K^0 \rightarrow \pi^+ + \pi^-. \qquad (5)$$

In the picture the dashed line $AB$ shows the path of the $K^0$ (continued past the vertex $B$ where it decays). This path is, of course, invisible in the photograph and has to be drawn by the student, by joining the two vertices.

The curvature of the charged tracks in Fig. 1 is due to a known magnetic field and lets the students measure the momenta of the charged particles. Using these measurements, the students can then carry out various interesting exercises: They can make a partial check that momentum is conserved in the $K^0$ decay, by verifying that the total momentum of the two pions has zero component perpendicular to the $K^0$ path; and they can calculate the $K^0$ mass, using the masses and measured momenta of the pions. Both of these exercises use measurements of the two angles shown as $\theta$ and $\phi$ in Fig. 1.

Under typical conditions, the measurements of $\theta$ and $\phi$ are strongly correlated. The uncertainty in $\theta$ and $\phi$ comes mainly from the uncertainty in the direction of the $K^0$ path, which must be found by joining the vertices $A$ and $B$. In a typical picture the distance between the vertices is of order 1 cm, while the tracks have a lateral blurring of about 1
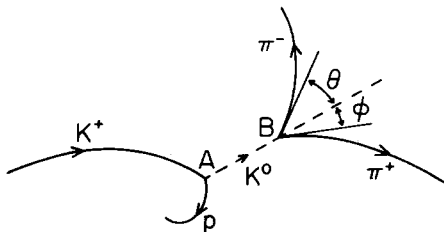


Fig. 1. The production of a $K^0$ meson at $A$ and its subsequent decay into two pions at $B$. The angles $\theta$ and $\phi$ are measured from the direction of the $K^0$ to the initial directions of the two pions.

mm. Under these conditions the direction of the $K^0$ track is uncertain by two or three degrees. On the other hand, the initial directions of the two pions are reasonably clearly defined, so that the total angle between them, $\theta + \phi$, suffers from very little uncertainty. Thus any underestimate of $\theta$ is accompanied by an equal overestimate of $\phi$ and vice versa; that is, $\theta$ and $\phi$ have a strong negative correlation.

To simplify our discussion, let us assume that the uncertainty in the direction of the $K^0$ track is the only source of uncertainty in the experiment. This is actually a good approximation: As I have already argued, the initial direction of the pions can be accurately determined, and the only other measurements are the pions' momenta, which can be found with very little uncertainty. With this assumption the correlation between $\theta$ and $\phi$ is complete; that is,

$$\sigma_{\theta\phi} = -\sigma_\theta^2 = -\sigma_\phi^2 \qquad (6)$$

and the correlation coefficient, $r = \sigma_{\theta\phi}/\sigma_\theta\sigma_\phi$, is $-1$.

In the light of this correlation, let us examine the two exercises mentioned above. First, conservation of momentum in the $K^0$ decay requires that the total transverse momentum $p_t$ of the final pions should be zero:

$$p_t = p_+ \sin\theta - p_- \sin\phi = 0, \qquad (7)$$

where $p_+$ and $p_-$ denote the momenta of the $\pi^+$ and $\pi^-$. Naturally the students' measured values for $p_t$ are not exactly zero, and the question that they must ask is this: Is the measured value of $p_t$ consistent with the expected value $p_t = 0$, given the uncertainty $\sigma_{p_t}$ in $p_t$?

A correct evaluation of $\sigma_{p_t}$ uses (2) to give

$$\sigma_{p_t}^2 = \left(\frac{\partial p_t}{\partial\theta}\right)^2 \sigma_\theta^2 + 2\frac{\partial p_t}{\partial\theta}\frac{\partial p_t}{\partial\phi}\sigma_{\theta\phi} + \left(\frac{\partial p_t}{\partial\phi}\right)^2 \sigma_\phi^2$$

or, using (6),

$$\sigma_{p_t}^2 = \left(\frac{\partial p_t}{\partial\theta} - \frac{\partial p_t}{\partial\phi}\right)^2 \sigma_\theta^2$$

$$= (p_+ \cos\theta + p_- \cos\phi)^2 \sigma_\theta^2.$$

In all of the events used in our laboratory, both $\theta$ and $\phi$ are less than 90°, so that both $p_+ \cos\theta$ and $p_- \cos\phi$ are positive. Thus

$$\sigma_{p_t} = (p_+ \cos\theta + p_- \cos\phi)\sigma_\theta. \qquad (8)$$

This shows that the uncertainty in $p_t$ is equal to its maximum possible value, as given by the linear sum (3).

If the students estimate the uncertainty $\sigma_{p_t}$ using the quadratic sum (1) they will get an answer which is too small by a factor of order $\sqrt{2}$. In my section of this laboratory none of the students noticed that $\theta$ and $\phi$ are correlated (nor did I until much later[4]) and several of the better students found that their values of $p_t$ were somewhat larger than their estimated value of $\sigma_{p_t}$. It is presumably fair to attribute at least part of this trend to their having overlooked $\sigma_{\theta\phi}$ and their consequent underestimation of $\sigma_{p_t}$.

The effect of the correlation in $\theta$ and $\phi$ is more dramatic in the determination of the $K^0$ mass. This is calculated from the relativistic relation

$$(mc^2)^2 = E^2 - (pc)^2, \qquad (9)$$

where $E$ and $p$ are the $K^0$ energy and momentum, which are found by assuming conservation of energy and momentum in the $K^0$ decay. That is, $E$ is the sum of the pions' energies,

$$E = E_+ + E_- \qquad (10)$$

while $p$ is the total longitudinal momentum of the two pions,

$$p = p_+ \cos\theta + p_- \cos\phi. \tag{11}$$

We are assuming that the pion momenta $p_+$ and $p_-$ can be measured with negligible uncertainty. The same is therefore true of the pion energies, and the only uncertainty in the $K^0$ mass $m$ comes from the angles $\theta$ and $\phi$ in Eq. (11) for the $K^0$ momentum $p$.

Given that $\theta$ and $\phi$ are correlated we must find the uncertainty in $p$ using Eq. (2), which gives

$$\sigma_p^2 = \left(\frac{\partial p}{\partial\theta}\right)^2\sigma_\theta^2 + 2\frac{\partial p}{\partial\theta}\frac{\partial p}{\partial\phi}\sigma_{\theta\phi} + \left(\frac{\partial p}{\partial\phi}\right)^2\sigma_\phi^2$$

or, since $\sigma_{\theta\phi} = -\sigma_\theta^2 = -\sigma_\phi^2$,

$$\sigma_p^2 = (p_+ \sin\theta - p_- \sin\phi)^2\sigma_\theta^2. \tag{12}$$

The term in parentheses will be recognized as $p_t$, the transverse momentum of the two pions, which, if we assume conservation of momentum, is zero.

This result does not, of course, mean that $\sigma_p$ is really zero. In the first place, all of our formulae are correct only to lowest order in $\sigma_\theta$. Thus the fact that the right-hand side of (12) is zero implies only that $\sigma_p$ is zero to lowest order in $\sigma_\theta$, and that we should probably not have neglected the smaller uncertainties in the other measurements. As a practical matter, the important conclusion is that $\sigma_p$ is *much smaller* than one would expect on the basis of naive use of the quadratic sum to estimate $\sigma_p$. In my section of the laboratory, several students found values for the $K^0$ mass that were appreciably better than they expected. Once again, it is presumably fair to attribute this to their overlooking the correlation of $\theta$ and $\phi$ and their consequent *overestimation* of the uncertainty concerned.[5]

## III. CORRELATION OF LEAST-SQUARES COEFFICIENTS

In my second example, the two correlated variables are not themselves measured directly. Instead they are the least-squares estimates for the coefficients $A$ and $B$ in a linear relation

$$y = A + Bx \tag{13}$$

based on $N$ measured pairs $(x_i, y_i)$, with $i = 1,..., N$. In the simplest case, which I consider here, the measured numbers $x_1,..., x_N$ have negligible uncertainty, while $y_1,..., y_N$ are all equally uncertain, with common standard deviation $\sigma_y$. In this case, the best estimates for $A$ and $B$ are well known to be[6]

$$A = [(\Sigma x_i^2)(\Sigma y_i) - (\Sigma x_i)(\Sigma x_i y_i)]/\Delta \tag{14}$$

and

$$B = [N(\Sigma x_i y_i) - (\Sigma x_i)(\Sigma y_i)]/\Delta, \tag{15}$$

where

$$\Delta = N(\Sigma x_i^2) - (\Sigma x_i)^2$$
$$= N\Sigma(x_i - \bar{x})^2. \tag{16}$$

The uncertainties in these estimates are given by

$$\sigma_A^2 = \sigma_y^2(\Sigma x_i^2)/\Delta \tag{17}$$

and

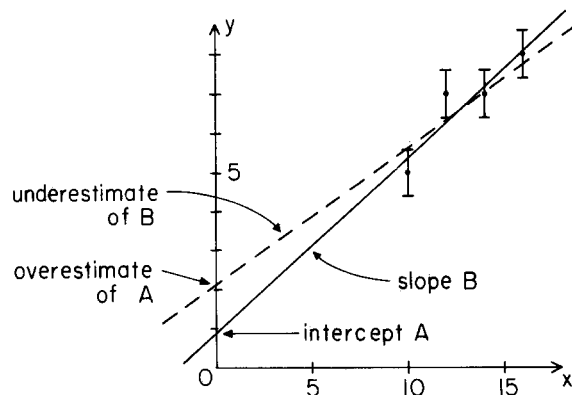$$\sigma_B^2 = \sigma_y^2 N/\Delta. \tag{18}$$



Fig. 2. Four measured points $(x_i, y_i)$ and the least-squares best line $y = A + Bx$ (solid line). The error bars show the uncertainties $\pm \sigma_y$ in the measured values $y_1,..., y_4$. The dashed line shows how an underestimate of the slope $B$ entails an overestimate of the intercept $A$.

In a typical application the estimates (14) and (15) are used to predict a value of $y$ for some chosen value of $x$. I shall denote this predicted value by

$$\hat{y} = A + Bx. \tag{19}$$

Note that the caret distinguishes the predicted value, $\hat{y}$, as defined by (19). Here, and from now on, $A$ and $B$ denote the estimates (14) and (15) [as opposed to the hypothetical "true" values in (13)].

The uncertainty in the predicted value $\hat{y}$ can be found in several ways: Perhaps the most straightforward is to substitute (14) and (15) into (19) to express $\hat{y}$ in terms of the measured quantities $y_1,..., y_N$. The uncertainty in $\hat{y}$ can then be found by error propagation—formula (1)—from the uncertainties in the uncorrelated measurements $y_1,..., y_N$.[7] Perhaps the neatest method is due to Birge,[8] who pointed out that for the special case that $x = 0$, $\sigma_{\hat{y}}$ is just $\sigma_A$; by making a suitable translation of origin, one can therefore find $\sigma_{\hat{y}}$ for any value of $x$.

The method of finding $\sigma_{\hat{y}}$ that I wish to discuss here is this: Having once found the two coefficients $A$ and $B$ with their uncertainties, it is natural to regard $A$ and $B$ as given, known quantities, and to calculate $\sigma_{\hat{y}}$ directly from (19) in terms of the uncertainties in $A$ and $B$. However, the two quantities $A$ and $B$ are correlated, as we now discuss, and one needs to know the covariance $\sigma_{AB}$ as well as $\sigma_A$ and $\sigma_B$ before one can find $\sigma_{\hat{y}}$ correctly.[9]

It is easy to understand why $A$ and $B$ are correlated. Figure 2 shows a typical set of data and the least-squares best fit (solid line). It is quite evident that, with these data, an *underestimate* of the slope $B$ (dashed line) will entail an *overestimate* of the intercept $A$, and vice versa. That is, the quantities $A$ and $B$ have (in this case) a strong negative correlation.

The covariance $\sigma_{AB}$ is easily evaluated. We imagine making a whole series of measurements of the $N$ variables $y_1,..., y_N$ (always for the same set $x_1,..., x_N$). Given that the measurements of $y_1,..., y_N$ are themselves uncorrelated,

$$\sigma_{AB} = \sum_{i=1}^{N}\frac{\partial A}{\partial y_i}\frac{\partial B}{\partial y_i}\sigma_{y_i}^2.$$

(This formula corresponds to the better-known formula for the variance,

$$\sigma_A^2 = \sum_{i=1}^{N}\left(\frac{\partial A}{\partial y_i}\right)^2\sigma_{y_i}^2,$$

and is derived in exactly the same way—see, for example, Ref. 1, pp. 176–177.) Substituting (14) and (15) for $A$ and $B$ and putting all of the $\sigma_{y_i}$ equal to $\sigma_y$, we find, after a little algebra,

$$\sigma_{AB} = -\sigma_y^2(\Sigma x_i)/\Delta. \qquad (20)$$

Knowing $\sigma_A$, $\sigma_B$, and $\sigma_{AB}$, we can now calculate $\sigma_{\hat{y}}$ using formula (2). After some more algebra, this gives

$$\sigma_{\hat{y}}^2 = \sigma_A^2 + 2x\sigma_{AB} + x^2\sigma_B^2$$

$$= \sigma_y^2 \sum (x - x_i)^2/\Delta \qquad (21)$$

$$= \frac{\sigma_y^2}{N} \frac{\Sigma (x - x_i)^2}{\Sigma (\bar{x} - x_i)^2}. \qquad (22)$$

These answers have several interesting features. First, if we choose a value of $x$ far away from all of the measured points, then

$$\Sigma (x - x_i)^2 \approx \Sigma (x - \bar{x})^2 = N(x - \bar{x})^2$$

and (21) gives

$$\sigma_{\hat{y}}^2 \approx \sigma_y^2 N(x - \bar{x})^2/\Delta = \sigma_B^2 (x - \bar{x})^2.$$

Therefore

$$\sigma_{\hat{y}} \approx \sigma_B |x - \bar{x}|.$$

This is just what we should expect: When $x$ is far away from the measured points we are extrapolating the line in Fig. 2 a long way. In this case the dominant uncertainty is the uncertainty in the slope $B$, and this contributes an uncertainty $\sigma_B |x - \bar{x}|$ to the predicted value $\hat{y}$.

More interesting, for our present purposes, is the case that we wish to predict $y$ for a value of $x$ close to the measured points $x_1,..., x_N$. Whatever the value of $x$, the right-hand fraction in (22) is never less than 1:

$$\frac{\Sigma (x - x_i)^2}{\Sigma (\bar{x} - x_i)^2} \geqslant 1. \qquad (23)$$

Thus (22) implies that

$$\sigma_{\hat{y}} \geqslant \sigma_y/\sqrt{N}, \qquad (24)$$

a result one might have anticipated. (The prediction $\hat{y}$ is based on $N$ measurements of $y$, so cannot possibly be more certain than $\sigma_y/\sqrt{N}$.) The fraction (23) is *equal* to 1 if (and only if) $x = \bar{x}$, and $\sigma_{\hat{y}}$ is therefore minimum at $x = \bar{x}$. That is, the most favorable place to predict $y$ is the centroid, $x = \bar{x}$, of the original measurements, again as one might have guessed.

Finally, if we were to overlook the correlation of $A$ and $B$, formula (1) would give

$$\sigma_{\hat{y}}^2 = \sigma_A^2 + x^2\sigma_B^2 \qquad (25)$$

$$= \frac{\sigma_y^2}{N} \frac{\Sigma (x^2 + x_i^2)}{\Sigma (\bar{x} - x_i)^2}. \qquad (26)$$

Comparison of this (incorrect) answer with (22) shows that by forgetting the correlation we have dropped the cross terms $-2xx_i$ from the sum in the numerator of (22). It is easy to see that one can construct examples in which the effect of this omission is very large. For example, let us consider the data of Fig. 2 and predict the value of $y$ at $x = \bar{x} = 13$ (the most favorable place to predict $y$). In this case the correct uncertainty (22) is

$$\sigma_{\hat{y}} = \sigma_y/\sqrt{4} = 0.3,$$

an answer which a glance at the graph shows to be reasonable. If we were to forget the correlation and use (26) we would get the absurd answer

$$\sigma_{\hat{y}} = 2.4.$$

Evidently the correlation of $A$ and $B$ is an important effect in this instance.

A simple concrete example of these considerations is given by an experiment on simple harmonic motion in our elementary physics laboratory at Colorado. In this laboratory, students first find the spring constant $k$ of a spring by loading it with four or five different masses and plotting the length $y$ of the spring against the mass $m$ of the load. Since the force of the spring is $k(y - y_0)$, where $y_0$ is its unstretched length, it follows that $mg = k(y - y_0)$ or

$$y = y_0 + (g/k)m. \qquad (27)$$

Therefore, the graph of $y$ against $m$ should be a straight line, very similar to that of Fig. 2, with intercept equal to the unstretched length $y_0$ and slope $g/k$.

In this experiment the data satisfy the assumptions of this section almost perfectly. The measurements of the length $y$ are all equally uncertain, by a millimeter or two, while the masses $m$ are standard masses and so have very small uncertainties. Thus a natural way to find the spring constant and the unstretched length is to make a least-squares fit to the line (27).[10] This gives the slope, and hence $k$, with a small uncertainty; but, because of the need to extrapolate back to zero load, the unstretched length $y_0$ has a larger uncertainty than the original measurements of $y$.[11]

Having found the spring constant $k$, students pick one of their masses $m$ and study its oscillations on the end of the spring. For this purpose they need to know the length $y$ of the spring when loaded with the mass $m$. This length has already been measured directly, but can be found more accurately from their fit to the line (27), by calculating $\hat{y}$ for the chosen mass $m$. For our present purposes, the important point about this calculation is this: The final uncertainty $\sigma_{\hat{y}}$, as given correctly by (22), is *smaller* than the uncertainty $\sigma_y$ in the original measurements. But an incorrect evaluation of $\sigma_{\hat{y}}$ based on (25) (i.e., neglecting the correlation of $A$ and $B$) would give the absurd answer that $\sigma_{\hat{y}}$ is larger than $\sigma_y$.

## ACKNOWLEDGMENTS

[1]See, for example, J. R. Taylor, *An Introduction to Error Analysis* (University Science Books, Mill Valley, CA, 1982), p. 177.

[2]See, for example, S. L. Meyer, *Data Analysis* (Wiley, New York, 1975), p. 367; or P. R. Bevington, *Data Reduction and Error Analysis* (McGraw-Hill, New York, 1969), p. 161.

[3]For example, D. C. Baird, *Experimentation* (Prentice-Hall, Englewood Cliffs, NJ, 1962); N. C. Barford, *Experimental Measurements* (Addison-Wesley, Reading, MA, 1967); Y. Beers, *Theory of Error* (Addison-Wesley, Reading, MA, 1957); E. M. Pugh and G. H. Winslow, *The Analysis of Physical Measurements* (Addison-Wesley, Reading, MA, 1966); J. Topping, *Errors of Observation* (Chapman and Hall, London, 1962); H. D. Young, *Statistical Treatment of Experimental Data* (McGraw-Hill, New York, 1962).

[4] I am indebted to my colleague Bill Ford for drawing my attention to this correlation.

[5] There is another way to calculate the $K^0$ momentum which bypasses the whole problem of correlation between $\theta$ and $\phi$. If one assumes conservation of all components of momentum, then $\mathbf{p} = \mathbf{p}_+ + \mathbf{p}_-$ and hence $p^2 = p_+^2 + 2p_+p_- \cos(\theta + \phi) + p_-^2$. Since this involves only $\theta + \phi$ (not $\theta$ or $\phi$ separately), the large uncertainties in $\theta$ and $\phi$ should have no impact on the value of $p$, just as long as $\theta + \phi$ can be accurately measured. The instructions for this experiment have recently been rewritten to suggest finding $p$ this way.

[6] See, for example, Ref. 1, pp. 156–159.

[7] See, for example, Beers, Ref. 3, p. 43.

[8] R. T. Birge, Phys. Rev. **40**, 207 (1932), see p. 226.

[9] This approach (regarding $A$ and $B$ as given) is encouraged by the many texts (including my own) that list formulas for $A$, $B$, $\sigma_A$, and $\sigma_B$. The only text that I have found with a clear warning that one needs to remember the correlation $\sigma_{AB}$ is Meyer, Ref. 2, p. 367: "It is common, if sometimes inadequate, to quote $A \pm \sigma_A$ and $B \pm \sigma_B$ as 'results,' but we should keep the existence of *correlations* in mind."

[10] It should be admitted that in our elementary laboratory the students do not actually use the method of least squares. Nevertheless, this is an almost perfect simple application of the method.

[11] The unstretched length $y_0$ cannot be measured directly because of the need to open up the spring's coils, which are in contact when the spring is unloaded. For the same reason, one cannot use small loads and must therefore extrapolate an appreciable distance back to zero load.

# Energy balance in the superposition of light waves with lossless beam splitters

F. Pi and G. Orriols

*Departament de Física Fonamental, Universitat Autònoma de Barcelona, Bellaterra, Barcelona, Spain*

The energy balance in the superposition of light waves with a lossless beam splitter is analyzed along with its connection with some general relations between the parameters describing beam-splitter effects on the two incoming waves. These general relations are derived either from the electromagnetic theory or from energy conservation considerations. The apparent paradox of a symmetric beam splitter producing complementary outputs is explained.

## I. INTRODUCTION

When dealing with light interference phenomena the student is often confused with questions related to energy balance. Surprisingly enough, these kinds of questions are not usually dealt with in textbooks.

The most typical example of a striking situation appears when considering the superposition of two copropagating plane waves and the resultant irradiance deviates from the sum of the component irradiances. As a matter of fact this situation can easily be achieved in the laboratory by using, for instance, a Michelson-type interferometer with a collimated light beam as input. The experiment is particularly useful for the verification of the fundamental formula

$$I = I_1 + I_2 + 2\sqrt{I_1 I_2}\, \cos \psi$$

giving the irradiance of the superposed beams, since the irradiances $I_1$ and $I_2$ of the two beams can easily be determined by alternative screening of the interferometer arms, and the phase difference $\psi$ of the waves varied by displacing the movable mirror. However, at first sight, the variation of the output beam power with $\psi$ poses an energy balance problem that must be explained to the students.

As shown in Fig. 1, the first point to be noticed is the existence of a second light output returning to the source from the interferometer. This second output can also be detected and any hypothetical influence on the source avoided if, for instance, the input beam direction is somewhat skew with respect to the interferometer horizontal plane.[1] In such a way the relative behavior of the two interference outputs can be analyzed as a function of the optical path difference between the interferometer arms. It is a frequent idea that the two outputs are complementary to one another. Nevertheless, if the experiment is carried out such an expectation will probably not prove to be so. The relative behavior of the outputs depends on the kind of beam splitter used in the interferometer. Complementary outputs are certainly observed with any lossless beam splitter, but not when energy is lost in the beam-splitter coating.

It is the aim of this paper to analyze the energy balance in the superposition of light waves with a lossless beam splitter and its connection with some general relations between the parameters describing beam-splitter effects on the two incoming waves. These general properties, satisfied for incident angles producing complete superposition at the outputs, are derived either from energy conservation considerations (Sec. III) or from the electromagnetic theory (Sec. IV). The particular case of a symmetric beam splitter is considered in some detail (Sec. V), since a symmetric device yielding complementary outputs may be somewhat surprising.

Early work on the subject deals with multiple-beam interference from a semireflecting film. As earlier as 1906 Hamy[2] established the relation that must exist between phase changes on transmission and reflection in order to obtain complementary patterns in transmitted or reflected light. It is this same relation that is required for complementary outputs in the beam-splitter interferometer. The