# III

# *Foundations of geometrical optics*

### 3.1 Approximation for very short wavelengths

THE electromagnetic field associated with the propagation of visible light is character-ized by very rapid oscillations (frequencies of the order of $10^{14}$ s$^{-1}$) or, what amounts to the same thing, by the smallness of the wavelength (of order $10^{-5}$ cm). It may therefore be expected that a good first approximation to the propagation laws in such cases may be obtained by a complete neglect of the finiteness of the wavelength. It is found that for many optical problems such a procedure is entirely adequate; in fact, phenomena which can be attributed to departures from this approximate theory (so-called diffraction phenomena, studied in Chapter VIII) can only be demonstrated by means of carefully conducted experiments.

The branch of optics which is characterized by the neglect of the wavelength, i.e. that corresponding to the limiting case $\lambda_0 \to 0$, is known as *geometrical optics*,[*] since in this approximation the optical laws may be formulated in the language of geometry. The energy may then be regarded as being transported along certain curves (light rays). A physical model of a pencil of rays may be obtained by allowing the light from a source of negligible extension to pass through a very small opening in an opaque screen. The light which reaches the space behind the screen will fill a region the boundary of which (the edge of the pencil) will, at first sight, appear to be sharp. A more careful examination will reveal, however, that the light intensity near the boundary varies rapidly but continuously from darkness in the shadow to lightness in the illuminated region, and that the variation is not monotonic but is of an oscillatory character, manifested by the appearance of bright and dark bands, called diffraction fringes. The region in which this rapid variation takes place is only of the order of magnitude of the wavelength. Hence, as long as this magnitude is neglected in comparison with the dimensions of the opening, we may speak of a sharply bounded pencil of rays.[†] On reducing the size of the opening down to the dimensions of the

---

[*] The historical development of geometrical optics is described by M. Herzberger, *Strahlenoptik* (Berlin, Springer, 1931), p. 179; *Z. Instrumentenkunde*, **52** (1932), 429–435, 485–493, 534–542, C. Carathéodory, *Geometrische Optik* (Berlin, Springer, 1937) and E. Mach, *The Principles of Physical Optics, A Historical and Philosophical Treatment* (First German edition 1913, English translation: London, Methuen, 1926; reprinted by Dover Publications, New York, 1953).

[†] That the boundary becomes sharp in the limit as $\lambda_0 \to 0$ was first shown by G. Kirchhoff, *Vorlesungen ü. Math. Phys.*, Vol. 2 (*Mathematische Optik*) (Leipzig, Tuebner, 1891), p. 33. See also B. B. Baker and E. T. Copson, *The Mathematical Theory of Huygens' Principle* (Oxford, Clarendon Press, 2nd edition, 1950), p. 79, and A. Sommerfeld, *Optics* (New York, Academic Press, 1954), §35.

116

wavelength phenomena appear which need more refined study. If, however, one considers only the limiting case of negligible wavelengths, no restriction on the size of the opening is imposed, and we may say that an opening of vanishingly small dimensions defines an infinitely thin pencil — the light ray. It will be shown that the variation in the cross-section of a pencil of rays is a measure of the variation of the intensity of the light. Moreover, it will be seen that it is possible to associate a state of polarization with each ray, and to study its variation along the ray.

Further it will be seen that for small wavelengths the field has the same general character as that of a plane wave, and, moreover, that within the approximation of geometrical optics the laws of refraction and reflection established for plane waves incident upon a plane boundary remain valid under more general conditions. Hence if a light ray falls on a sharp boundary (e.g. the surface of a lens) it is split into a reflected ray and a transmitted ray, and the changes in the state of polarization as well as the reflectivity and transmissivity may be calculated from the corresponding formulae for plane waves.

The preceding remarks imply that, when the wavelength is small enough, the sum total of optical phenomena may be deduced from geometrical considerations, by determining the paths of the light rays and calculating the associated intensity and polarization. We shall now formulate the appropriate laws by considering the implications of Maxwell's equations when $\lambda_0 \to 0$.[*]

### 3.1.1 Derivation of the eikonal equation

We consider a general time-harmonic field

$$\left.\begin{array}{l} \mathbf{E}(\mathbf{r},\, t) = \mathbf{E}_0(\mathbf{r})\mathrm{e}^{-\mathrm{i}\omega t}, \\[4pt] \mathbf{H}(\mathbf{r},\, t) = \mathbf{H}_0(\mathbf{r})\mathrm{e}^{-\mathrm{i}\omega t}, \end{array}\right\} \tag{1}$$

in a nonconducting isotropic medium. $\mathbf{E}_0$ and $\mathbf{H}_0$ denote complex vector functions of positions, and, as explained in §1.4.3, the real parts of the expressions on the right-hand side of (1) are understood to represent the fields.

The complex vectors $\mathbf{E}_0$ and $\mathbf{H}_0$ will satisfy Maxwell's equations in their time-free form, obtained on substituting (1) into §1.1 (1)–(4). In regions free of currents and charges ($\mathbf{j} = \rho = 0$), these equations are

$$\operatorname{curl} \mathbf{H}_0 + \mathrm{i}k_0\varepsilon\mathbf{E}_0 = 0, \tag{2}$$

$$\operatorname{curl} \mathbf{E}_0 - \mathrm{i}k_0\mu\mathbf{H}_0 = 0, \tag{3}$$

$$\operatorname{div} \varepsilon\mathbf{E}_0 = 0, \tag{4}$$

$$\operatorname{div} \mu\mathbf{H}_0 = 0. \tag{5}$$

---

[*] It was first shown by A. Sommerfeld and J. Runge, *Ann. d. Physik*, **35** (1911), 289, using a suggestion of P. Debye, that the basic equation of geometrical optics [the eikonal equation (15b)] may be derived from the (scalar) wave equation in the limiting case $\lambda_0 \to 0$. Generalizations which take into account the vectorial character of the electromagnetic field are due to W. S. Ignatowsky, *Trans. State Opt. Institute* (*Petrograd*), **1** (1919), No 3, 30; V. A. Fock, *ibid.*, **3** (1924), 3; S. M. Rytov, *Compt. Rend. (Doklady) Acad. Sci. URSS*, **18** (1938), 263; N. Arley, *Det. Kgl. Danske Videns Selsk.*, **22** (1945), No. 8; F. G. Friedlander, *Proc. Cambr. Phil. Soc.*, **43** (1947), 284; K. Suchy, *Ann. d. Physik*, **11** (1952), 113, *ibid.*, **12** (1953), 423, and *ibid.*, **13** (1953), 178; R. S. Ingarden and A. Krzywicki, *Acta Phys. Polonica*, **14** (1955), 255.

Here the material relations $\mathbf{D} = \varepsilon\mathbf{E}$, $\mathbf{B} = \mu\mathbf{H}$ have been used and, as before, $k_0 = \omega/c = 2\pi/\lambda_0$, $\lambda_0$ being the vacuum wavelength.

We have seen that a homogeneous plane wave in a medium of refractive index $n = \sqrt{\varepsilon\mu}$, propagated in the direction specified by the unit vector $\mathbf{s}$, is represented by

$$\mathbf{E}_0 = \mathbf{e}\mathrm{e}^{\mathrm{i}k_0 n(\mathbf{s}\cdot\mathbf{r})}, \qquad \mathbf{H}_0 = \mathbf{h}\mathrm{e}^{\mathrm{i}k_0 n(\mathbf{s}\cdot\mathbf{r})}, \tag{6}$$

where $\mathbf{e}$ and $\mathbf{h}$ are constant, generally complex vectors. For a (monochromatic) electric dipole field in the vacuum we found (see §2.2) that

$$\mathbf{E}_0 = \mathbf{e}\mathrm{e}^{\mathrm{i}k_0 r}, \qquad \mathbf{H}_0 = \mathbf{h}\mathrm{e}^{\mathrm{i}k_0 r}, \tag{7}$$

$r$ being the distance from the dipole. Here $\mathbf{e}$ and $\mathbf{h}$ are no longer constant vectors, but at distances sufficiently far away from the dipole ($r \gg \lambda_0$) these vectors are, with suitable normalization of the dipole moment, independent of $k_0$.

These examples suggest that in regions which are many wavelengths away from the sources we may represent more general types of fields in the form

$$\mathbf{E}_0 = \mathbf{e}(\mathbf{r})\mathrm{e}^{\mathrm{i}k_0\mathcal{S}(\mathbf{r})}, \qquad \mathbf{H}_0 = \mathbf{h}(\mathbf{r})\mathrm{e}^{\mathrm{i}k_0\mathcal{S}(\mathbf{r})}, \tag{8}$$

where $\mathcal{S}(\mathbf{r})$, 'the optical path', is a real scalar function of position, and $\mathbf{e}(\mathbf{r})$ and $\mathbf{h}(\mathbf{r})$ are vector functions of position, which may in general be complex.[*] With (8) as a trial solution, Maxwell's equations lead to a set of relations between $\mathbf{e}$, $\mathbf{h}$ and $\mathcal{S}$. It will be shown that for large $k_0$ (small $\lambda_0$) these relations demand that $\mathcal{S}$ should satisfy a certain differential equation, which is independent of the amplitude vectors $\mathbf{e}$ and $\mathbf{h}$.

From (8), using well-known vector identities,

$$\mathrm{curl}\,\mathbf{H}_0 = (\mathrm{curl}\,\mathbf{h} + \mathrm{i}k_0\,\mathrm{grad}\,\mathcal{S} \times \mathbf{h})\mathrm{e}^{\mathrm{i}k_0\mathcal{S}}, \tag{9}$$

$$\mathrm{div}\,\mu\mathbf{H}_0 = (\mu\,\mathrm{div}\,\mathbf{h} + \mathbf{h}\cdot\mathrm{grad}\,\mu + \mathrm{i}k_0\mu\mathbf{h}\cdot\mathrm{grad}\,\mathcal{S})\mathrm{e}^{\mathrm{i}k_0\mathcal{S}}, \tag{10}$$

with similar expressions for $\mathrm{curl}\,\mathbf{E}_0$ and $\mathrm{div}\,\varepsilon\mathbf{E}_0$. Hence (2)–(5) become

$$\mathrm{grad}\,\mathcal{S} \times \mathbf{h} + \varepsilon\mathbf{e} = -\frac{1}{\mathrm{i}k_0}\,\mathrm{curl}\,\mathbf{h}, \tag{11}$$

$$\mathrm{grad}\,\mathcal{S} \times \mathbf{e} - \mu\mathbf{h} = -\frac{1}{\mathrm{i}k_0}\,\mathrm{curl}\,\mathbf{e}, \tag{12}$$

$$\mathbf{e}\cdot\mathrm{grad}\,\mathcal{S} = -\frac{1}{\mathrm{i}k_0}(\mathbf{e}\cdot\mathrm{grad}\ln\varepsilon + \mathrm{div}\,\mathbf{e}), \tag{13}$$

$$\mathbf{h}\cdot\mathrm{grad}\,\mathcal{S} = -\frac{1}{\mathrm{i}k_0}(\mathbf{h}\cdot\mathrm{grad}\ln\mu + \mathrm{div}\,\mathbf{h}). \tag{14}$$

We are interested in the solution for very large values of $k_0$. Hence as long as the multiplicative factors of $1/\mathrm{i}k_0$ on the right-hand side are not exceptionally large they may be neglected, and the equations then reduce to

---

[*] Complex $\mathbf{e}$ and $\mathbf{h}$ are necessary, if all possible states of polarization are to be included. According to §1.4 (75) real $\mathbf{e}$ and $\mathbf{h}$ correspond to fields which are linearly polarized.

$$\text{grad}\,\mathcal{S} \times \mathbf{h} + \varepsilon \mathbf{e} = 0, \tag{11a}$$

$$\text{grad}\,\mathcal{S} \times \mathbf{e} - \mu \mathbf{h} = 0, \tag{12a}$$

$$\mathbf{e} \cdot \text{grad}\,\mathcal{S} = 0, \tag{13a}$$

$$\mathbf{h} \cdot \text{grad}\,\mathcal{S} = 0. \tag{14a}$$

We can confine our attention to (11a) and (12a) alone, since (13a) and (14a) follow from them on scalar multiplication with grad $\mathcal{S}$. Now (11a) and (12a) may be regarded as a set of six simultaneous homogeneous linear scalar equations for the Cartesian components $e_x$, $h_x$, ... of $\mathbf{e}$ and $\mathbf{h}$. These simultaneous equations have nontrivial solutions only if a consistency condition (the vanishing of the associated determinant) is satisfied. This condition may be obtained simply by eliminating $\mathbf{e}$ or $\mathbf{h}$ between (11a) and (12a). Substituting for $\mathbf{h}$ from (12a), (11a) becomes

$$\frac{1}{\mu}[(\mathbf{e} \cdot \text{grad}\,\mathcal{S})\,\text{grad}\,\mathcal{S} - \mathbf{e}(\text{grad}\,\mathcal{S})^2] + \varepsilon \mathbf{e} = 0.$$

The first term vanishes on account of (13a), and the equation then reduces, since $\mathbf{e}$ does not vanish everywhere, to

$$(\text{grad}\,\mathcal{S})^2 = n^2, \tag{15a}$$

or, written explicitly,

$$\left(\frac{\partial \mathcal{S}}{\partial x}\right)^2 + \left(\frac{\partial \mathcal{S}}{\partial y}\right)^2 + \left(\frac{\partial \mathcal{S}}{\partial z}\right)^2 = n^2(x, y, z), \tag{15b}$$

where as before $n = \sqrt{\varepsilon\mu}$ denotes the refractive index. The function $\mathcal{S}$ is often called the *eikonal**\** and (15b) is known as the *eikonal equation*; it is the basic equation of geometrical optics.† The surfaces

$$\mathcal{S}(\mathbf{r}) = \text{constant}$$

may be called the *geometrical wave surfaces* or the *geometrical wave-fronts*.‡

The eikonal equation was derived here by using the first-order Maxwell's equations, but it may also be derived from the second-order wave equations for the electric or magnetic field vectors. To show this one substitutes from (1) and (8) into the wave equation §1.2 (5) and obtains, after a straightforward calculation,

$$\boldsymbol{K}(\mathbf{e},\, \mathcal{S},\, n) + \frac{1}{\mathrm{i}k_0}\boldsymbol{L}(\mathbf{e},\, \mathcal{S},\, n,\, \mu) + \frac{1}{(\mathrm{i}k_0)^2}\boldsymbol{M}(\mathbf{e},\, \varepsilon,\, \mu) = 0, \tag{16}$$

where

---

\* The term *eikonal* (from Greek $\varepsilon\iota\kappa\tilde{\omega}\nu$ = image) was introduced in 1895 by H. Bruns to describe certain related functions (see p. 142), but has come to be used in a wider sense.

† The eikonal equation may also be regarded as the equation of the characteristics of the wave equations §1.2 (5) and §1.2 (6) for $\mathbf{E}$ and $\mathbf{H}$, and describes the propagation of discontinuities of the solutions of these equations. In geometrical optics we are, however, not concerned with the propagation of discontinuities but with time-harmonic (or nearly time-harmonic) solutions. The formal equivalence of the two interpretations is demonstrated in Appendix VI.

The eikonal equation will also be recognized as the Hamilton–Jacobi equation of the variational problem $\delta \int n\,\mathrm{d}s = 0$, the optical counterpart of which goes back to Fermat (see §3.3.2 and Appendix I).

‡ In future we shall drop the adjective 'geometrical' when there is no risk of confusion.

$$\boldsymbol{K}(\mathbf{e},\, \mathcal{S},\, n) = [n^2 - (\operatorname{grad} \mathcal{S})^2]\mathbf{e},$$

$$\boldsymbol{L}(\mathbf{e},\, \mathcal{S},\, n,\, \mu) = [\operatorname{grad} \mathcal{S} \cdot \operatorname{grad} \ln \mu - \nabla^2 \mathcal{S}]\mathbf{e} - 2[\mathbf{e} \cdot \operatorname{grad} \ln n]\operatorname{grad} \mathcal{S}$$
$$- 2[\operatorname{grad} \mathcal{S} \cdot \operatorname{grad}]\mathbf{e},$$

$$\boldsymbol{M}(\mathbf{e},\, \varepsilon,\, \mu) = \operatorname{curl} \mathbf{e} \times \operatorname{grad} \ln \mu - \nabla^2 \mathbf{e} - \operatorname{grad}(\mathbf{e} \cdot \operatorname{grad} \ln \varepsilon).$$

The corresponding equation involving $\mathbf{h}$, obtained on substitution into the wave equation §1.2 (6) for $\mathbf{H}$ (or more simply by using the fact that Maxwell's equations remain unchanged when $\mathbf{E}$ and $\mathbf{H}$ and simultaneously $\varepsilon$ and $-\mu$ are interchanged), is

$$\boldsymbol{K}(\mathbf{h},\, \mathcal{S},\, n) + \frac{1}{\mathrm{i}k_0} \boldsymbol{L}(\mathbf{h},\, \mathcal{S},\, n,\, \varepsilon) + \frac{1}{(\mathrm{i}k_0)^2} \boldsymbol{M}(\mathbf{h},\, \mu,\, \varepsilon) = 0. \tag{17}$$

For sufficiently large $k_0$ the second and third terms may in general be neglected; then $\boldsymbol{K} = 0$, giving again the eikonal equation. It will be seen later that the terms in the first power of $1/\mathrm{i}k_0$ in (16) and (17) also possess a physical interpretation.

It may be shown that in many cases of importance the spatial parts $\mathbf{E}_0$ and $\mathbf{H}_0$ of the field vectors may be developed into asymptotic series of the form[*]

$$\mathbf{E}_0 = \mathrm{e}^{\mathrm{i}k_0 \mathcal{S}} \sum_{m \geqslant 0} \frac{\mathbf{e}^{(m)}}{(\mathrm{i}k_0)^m}, \qquad \mathbf{H}_0 = \mathrm{e}^{\mathrm{i}k_0 \mathcal{S}} \sum_{m \geqslant 0} \frac{\mathbf{h}^{(m)}}{(\mathrm{i}k_0)^m}, \tag{18}$$

where $\mathbf{e}^{(m)}$ and $\mathbf{h}^{(m)}$ are functions of position, and $\mathcal{S}$ is the same function as before.[†] Geometrical optics corresponds to the leading terms of these expansions.

### 3.1.2 The light rays and the intensity law of geometrical optics

From (8), and from §1.4 (54) and §1.4 (55), it follows that the time averages of the electric and magnetic energy densities $\langle w_e \rangle$ and $\langle w_m \rangle$ are given by

$$\langle w_e \rangle = \frac{\varepsilon}{16\pi} \mathbf{e} \cdot \mathbf{e}^{\star}, \qquad \langle w_m \rangle = \frac{\mu}{16\pi} \mathbf{h} \cdot \mathbf{h}^{\star}. \tag{19}$$

Substitution for $\mathbf{e}^{\star}$ from (11a) and for $\mathbf{h}$ from (12a) gives

$$\langle w_e \rangle = \langle w_m \rangle = \frac{1}{16\pi}[\mathbf{e},\, \mathbf{h}^{\star},\, \operatorname{grad} \mathcal{S}], \tag{20}$$

the square bracket denoting the scalar triple product. Hence, *within the accuracy of geometrical optics, the time-averaged electric and magnetic energy densities are equal*.

The time average of the Poynting vector is obtained from (8) and §1.4 (56):

$$\langle \mathbf{S} \rangle = \frac{c}{8\pi} \mathcal{R}(\mathbf{e} \times \mathbf{h}^{\star}).$$

---

[*] We assume here that only one geometrical wave-front passes through each point. In some cases, for example when reflection takes place at obstacles present in the medium, several wave-fronts may pass through each point. The resulting field is then represented by the addition of series of the above type.

[†] The theory of such asymptotic expansions has its origin chiefly in the work of R. K. Luneburg, *Propagation of Electromagnetic Waves* (mimeographed lecture notes, New York University, 1947–1948). See also M. Kline, *Comm. Pure and Appl. Math.*, **4** (1951), 225; *ibid.*, **8** (1955), 595 and W. Braunbek, *Z. Naturforsch.*, **6** (1951), 672. A comprehensive account of the theory is given in M. Kline and I. W. Kay, *Electromagnetic Theory and Geometrical Optics* (New York, Interscience Publishers, 1965).

Using (12a), we obtain

$$\langle \mathbf{S} \rangle = \frac{c}{8\pi\mu} \{ (\mathbf{e} \cdot \mathbf{e}^{\star}) \operatorname{grad} \mathcal{S} - (\mathbf{e} \cdot \operatorname{grad} \mathcal{S}) \mathbf{e}^{\star} \}.$$

The last term vanishes on account of (13a) so that we have, if use is made of the expression for $\langle w_e \rangle$ and of Maxwell's relation $\varepsilon\mu = n^2$,

$$\langle \mathbf{S} \rangle = \frac{2c}{n^2} \langle w_e \rangle \operatorname{grad} \mathcal{S}. \tag{21}$$

Since $\langle w_e \rangle = \langle w_m \rangle$, the term $2\langle w_e \rangle$ represents the time average $\langle w \rangle$ of the total energy density (i.e. $\langle w \rangle = \langle w_e \rangle + \langle w_m \rangle$). Also, on account of the eikonal equation, $(\operatorname{grad} \mathcal{S})/n$ is a unit vector ($\mathbf{s}$ say),

$$\mathbf{s} = \frac{\operatorname{grad} \mathcal{S}}{n} = \frac{\operatorname{grad} \mathcal{S}}{|\operatorname{grad} \mathcal{S}|}, \tag{22}$$

and (21) shows that $\mathbf{s}$ is in the direction of the average Poynting vector. If, as before, we set $c/n = v$, (21) becomes

$$\langle \mathbf{S} \rangle = v \langle w \rangle \mathbf{s}. \tag{23}$$

Hence *the average Poynting vector is in the direction of the normal to the geometrical wave-front, and its magnitude is equal to the product of the average energy density and the velocity $v = c/n$.* This result is analogous to the relation §1.4 (9) for plane waves, and shows that *within the accuracy of geometrical optics the average energy density is propagated with the velocity $v = c/n$.*

The *geometrical light rays* may now be defined as the orthogonal trajectories to the geometrical wave-fronts $\mathcal{S} = $ constant. We shall regard them as oriented curves whose direction coincides everywhere with the direction of the average Poynting vector.[*] If $\mathbf{r}(s)$ denotes the position vector of a point $P$ on a ray, considered as a function of the length of arc $s$ of the ray, then $d\mathbf{r}/ds = \mathbf{s}$, and the equation of the ray may be written as

$$n \frac{d\mathbf{r}}{ds} = \operatorname{grad} \mathcal{S}. \tag{24}$$

From (13a) and (14a) it is seen that *the electric and magnetic vectors are at every point orthogonal to the ray.*

The meaning of (24) may be made clearer from the following remarks. Consider two neighbouring wave-fronts $\mathcal{S} = $ constant and $\mathcal{S} + d\mathcal{S} = $ constant (Fig. 3.1). Then

$$\frac{d\mathcal{S}}{ds} = \frac{d\mathbf{r}}{ds} \cdot \operatorname{grad} \mathcal{S} = n. \tag{25}$$

Hence the distance $ds$ between points on the opposite ends of a normal cutting the two wave-fronts is inversely proportional to the refractive index, i.e. directly proportional to $v$.

The integral $\int_C n \, ds$ taken along a curve $C$ is known as the *optical length* of the

---

[*] This definition of light rays is appropriate for isotropic media only. We shall see later (Chapter XV) that in an anisotropic medium the direction of the wave-front normal does not, in general, coincide with the direction of the Poynting vector.
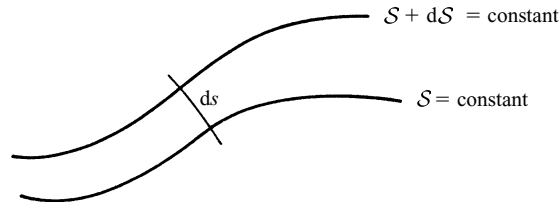
Fig. 3.1 Illustrating the meaning of the relation $n\mathbf{s} = \operatorname{grad}\mathcal{S}$.

curve. Denoting by square brackets the optical length of the *ray* which joins points $P_1$ and $P_2$, we have

$$[P_1 P_2] = \int_{P_1}^{P_2} n \, ds = \mathcal{S}(P_2) - \mathcal{S}(P_1). \tag{26}$$

Since, as we have seen, the average energy density is propagated with the velocity $v = c/n$ along the ray,

$$n \, ds = \frac{c}{v} \, ds = c \, dt,$$

where $dt$ is the time needed for the energy to travel the distance $ds$ along the ray; hence

$$[P_1 P_2] = c \int_{P_1}^{P_2} dt, \tag{27}$$

i.e. *the optical length $[P_1 P_2]$ is equal to the product of the vacuum velocity of light and the time needed for light to travel from $P_1$ to $P_2$.*

The intensity of light $I$ was defined as the absolute value of the time average of the Poynting vector. We therefore have from (23),

$$I = |\langle \mathbf{S} \rangle| = v \langle w \rangle, \tag{28}$$

and the conservation law §1.4 (57) gives

$$\operatorname{div}(I\mathbf{s}) = 0. \tag{29}$$

To see the implications of this relation we take a narrow tube formed by all the rays proceeding from an element $dS_1$ of a wave-front $\mathcal{S}(\mathbf{r}) = a_1$ ($a_1$ being a constant), and denote by $dS_2$ the corresponding element in which these rays intersect another wave-front $\mathcal{S}(\mathbf{r}) = a_2$ (Fig. 3.2). Integrating (29) throughout the tube and applying Gauss' theorem we obtain

$$\int I\mathbf{s} \cdot \boldsymbol{v} \, dS = 0, \tag{30}$$

$\boldsymbol{v}$ denoting the outward normal to the tube. Now

$$\mathbf{s} \cdot \boldsymbol{v} = \quad 1 \text{ on } dS_2,$$
$$= -1 \text{ on } dS_1,$$
$$= \quad 0 \text{ elsewhere,}$$

Fig. 3.2 Illustrating the intensity law of geometrical optics.

so that (30) reduces to

$$I_1 \, \mathrm{d}S_1 = I_2 \, \mathrm{d}S_2, \tag{31}$$

$I_1$ and $I_2$ denoting the intensity on $\mathrm{d}S_1$ and on $\mathrm{d}S_2$ respectively. Hence *IdS remains constant along a tube of rays*. This result expresses *the intensity law of geometrical optics*.

We shall see later that in a homogeneous medium the rays are straight lines. The intensity law may then be expressed in a somewhat different form. Assume first that $\mathrm{d}S_1$, and consequently also $\mathrm{d}S_2$, are bounded by segments of lines of curvature (see Fig. 3.3). If $R_1$ and $R_1'$ are the principal radii of curvature (§4.6.1) of the segments $A_1 B_1$ and $B_1 C_1$, then

$$A_1 B_1 = R_1 \, \mathrm{d}\theta, \qquad B_1 C_1 = R_1' \, \mathrm{d}\phi,$$

where $\mathrm{d}\theta$ and $\mathrm{d}\phi$ are the angles which $A_1 B_1$ and $B_1 C_1$ subtend at the respective centres of curvature $Q$ and $Q'$. Hence

$$\mathrm{d}S_1 = A_1 B_1 \cdot B_1 C_1 = R_1 R_1' \, \mathrm{d}\theta \, \mathrm{d}\phi; \tag{32}$$

similarly for an element $\mathrm{d}S_2$ in which the bundle of rays through $\mathrm{d}S_1$ meets another wave-front of the family,



$$QA_2 = R_2 = R_1 + l$$
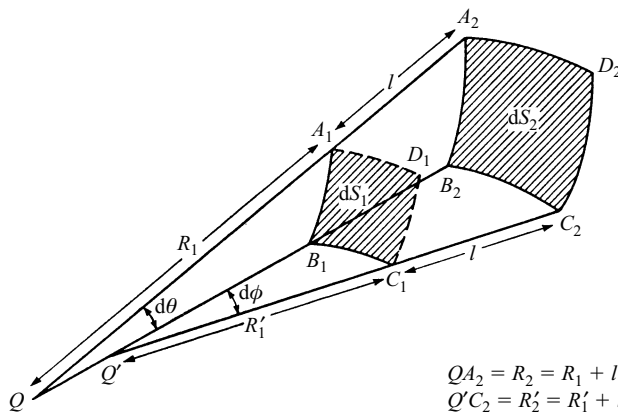$$Q'C_2 = R_2' = R_1' + l$$

Fig. 3.3 Illustrating the intensity law of geometrical optics for rectilinear rays.

$$dS_2 = A_2 B_2 \cdot B_2 C_2 = R_2 R_2' \, d\theta \, d\phi. \tag{33}$$

If $l$ is the distance between $dS_1$ and $dS_2$ measured along the rays, then

$$R_2 = R_1 + l, \qquad R_2' = R_1' + l,$$

and it follows that

$$\frac{I_2}{I_1} = \frac{dS_1}{dS_2} = \frac{R_1 R_1'}{R_2 R_2'} = \frac{R_1 R_1'}{(R_1 + l)(R_1' + l)}. \tag{34}$$

If the areas $dS_1$ and $dS_2$ are bounded by arbitrary curves, (34) still holds. This can immediately be seen if we regard them as made up of a number of elements bounded by lines of curvature, and sum their contributions.

If $R_1 \ll l$, $R_1' \ll l$, (34) reduces to

$$\frac{I_2}{I_1} = \frac{R_1 R_1'}{l^2}. \tag{35}$$

This formula is sometimes used in connection with the scattering of radiation.

The reciprocal $1/RR'$ of the product of the two principal radii of curvature is called the *Gaussian curvature* (or second curvature) of the surface. Eq. (34) shows that *the intensity at any point of a rectilinear ray is proportional to the Gaussian curvature of the wave-front which passes through that point.* In particular if all the (rectilinear) rays have a point in common, the wave-fronts are spheres centred at that point; then $R_1 = R_1'$, $R_2 = R_2'$ and we obtain (dropping the suffixes) the *inverse square law*

$$I = \frac{\text{constant}}{R^2}. \tag{36}$$

Returning to the general case of an arbitrary pencil of rays (curved or straight), we can write down an explicit expression in terms of the $\mathcal{S}$ function for the variation of the intensity along each ray. Substituting for $s$ from (22) into (29), and using the identities $\text{div}\, u\mathbf{v} = u \,\text{div}\, \mathbf{v} + \mathbf{v} \cdot \text{grad}\, u$, and $\text{div}\, \text{grad} = \nabla^2$, we obtain

$$\frac{I}{n} \nabla^2 \mathcal{S} + \text{grad}\, \mathcal{S} \cdot \text{grad}\, \frac{I}{n} = 0.$$

This may also be written as

$$\nabla^2 \mathcal{S} + \text{grad}\, \mathcal{S} \cdot \text{grad}\, \ln \frac{I}{n} = 0. \tag{37}$$

Let us now introduce the operator

$$\frac{\partial}{\partial \tau} = \text{grad}\, \mathcal{S} \cdot \text{grad}, \tag{38}$$

where $\tau$ is a parameter which specifies position along the ray. Then (37) may be written as

$$\frac{\partial}{\partial \tau} \ln \frac{I}{n} = -\nabla^2 \mathcal{S}$$

whence, on integration,

$$I = n e^{-\int^{\tau} \nabla^2 \mathcal{S} \, d\tau}.$$

But by (38), (15) and (25),

$$\mathrm{d}\tau = \frac{\mathrm{d}\mathcal{S}}{(\mathrm{grad}\,\mathcal{S})^2} = \frac{1}{n^2}\,\mathrm{d}\mathcal{S} = \frac{1}{n}\,\mathrm{d}s, \tag{39}$$

so that we finally obtain the following expressions for the ratio of the intensities at any two points of a ray:

$$\frac{I_2}{I_1} = \frac{n_2}{n_1}\,\mathrm{e}^{-\int_{\mathcal{S}_1}^{\mathcal{S}_2}\frac{\nabla^2 \mathcal{S}}{n^2}\,\mathrm{d}\mathcal{S}} = \frac{n_2}{n_1}\,\mathrm{e}^{-\int_{s_1}^{s_2}\frac{\nabla^2 \mathcal{S}}{n}\,\mathrm{d}s}, \tag{40}$$

the integrals being taken along the ray.[*]

### 3.1.3 Propagation of the amplitude vectors

We have seen that, when the wavelength is sufficiently small, the transport of energy may be represented by means of a simple hydrodynamical model which may be completely described in terms of the real scalar function $\mathcal{S}$, this function being a solution of the eikonal equation (15). According to traditional terminology, one understands by geometrical optics this approximate picture of energy propagation, using the concept of rays and wave-fronts. In other words polarization properties are excluded. The reason for this restriction is undoubtedly due to the fact that the simple laws of geometrical optics concerning rays and wave-fronts were known from experiments long before the electromagnetic theory of light was established. It is, however, possible, and from our point of view quite natural, to extend the meaning of geometrical optics to embrace also certain geometrical laws relating to the propagation of the 'amplitude vectors' **e** and **h**. These laws may be easily deduced from the wave equations (16)–(17).

Since $\mathcal{S}$ satisfies the eikonal equation, it follows that $\boldsymbol{K} = 0$, and we see that when $k_0$ is sufficiently large ($\lambda_0$ small enough), only the $\boldsymbol{L}$-terms need to be retained in (16) and (17). Hence, in the present approximation, the amplitude vectors and the eikonal are connected by the relations $\boldsymbol{L} = 0$. If we use again the operator $\partial/\partial\tau$ introduced by (38), the equations $\boldsymbol{L} = 0$ become

$$\frac{\partial \mathbf{e}}{\partial \tau} + \frac{1}{2}\left(\nabla^2 \mathcal{S} - \frac{\partial \ln \mu}{\partial \tau}\right)\mathbf{e} + (\mathbf{e}\cdot\mathrm{grad}\,\ln n)\,\mathrm{grad}\,\mathcal{S} = 0, \tag{41}$$

$$\frac{\partial \mathbf{h}}{\partial \tau} + \frac{1}{2}\left(\nabla^2 \mathcal{S} - \frac{\partial \ln \varepsilon}{\partial \tau}\right)\mathbf{h} + (\mathbf{h}\cdot\mathrm{grad}\,\ln n)\,\mathrm{grad}\,\mathcal{S} = 0. \tag{42}$$

These are the required *transport equations* for the variation of **e** and **h** along each ray. The implications of these equations can best be understood by examining separately the variation of the magnitude and of the direction of these vectors.

---

[*] It has been shown by M. Kline, *Comm. Pure and Appl. Maths.*, **14** (1961), 473 that the intensity ratio (40) may be expressed in terms of an integral which involves the principal radii of curvature of the associated wavefronts. Kline's formula is a natural generalization, to inhomogeneous media, of the formula (34). See also M. Kline and I. W. Kay, *ibid*, 184.

We multiply (41) scalarly by $\mathbf{e}^\star$ and add to the resulting equation the corresponding equation obtained by taking the complex conjugate. This gives

$$\frac{\partial}{\partial \tau}(\mathbf{e} \cdot \mathbf{e}^\star) + \left(\nabla^2 \mathcal{S} - \frac{\partial \ln \mu}{\partial \tau}\right)\mathbf{e} \cdot \mathbf{e}^\star = 0. \tag{43}$$

On account of the identity $\operatorname{div} u\mathbf{v} = u \operatorname{div} \mathbf{v} + \mathbf{v} \cdot \operatorname{grad} u$, the second and third term may be combined as follows:

$$\nabla^2 \mathcal{S} - \frac{\partial \ln \mu}{\partial \tau'} = \nabla^2 \mathcal{S} - \operatorname{grad} \mathcal{S} \cdot \operatorname{grad} \ln \mu = \mu \operatorname{div}\left(\frac{1}{\mu} \operatorname{grad} \mathcal{S}\right). \tag{44}$$

Integrating (43) along a ray, the following expression for the ratio of $\mathbf{e} \cdot \mathbf{e}^\star$ at any two points of the ray is obtained:*

$$\frac{(\mathbf{e} \cdot \mathbf{e}^\star)_2}{(\mathbf{e} \cdot \mathbf{e}^\star)_1} = \mathrm{e}^{-\int_{\tau_1}^{\tau_2} \mu \operatorname{div}\left(\frac{1}{\mu} \operatorname{grad} \mathcal{S}\right)\mathrm{d}\tau} = \mathrm{e}^{-\int_{s_1}^{s_2} \sqrt{\frac{\mu}{\varepsilon}} \operatorname{div}\left(\frac{1}{\mu} \operatorname{grad} \mathcal{S}\right)\mathrm{d}s}. \tag{45}$$

Similarly

$$\frac{(\mathbf{h} \cdot \mathbf{h}^\star)_2}{(\mathbf{h} \cdot \mathbf{h}^\star)_1} = \mathrm{e}^{-\int_{s_1}^{s_2} \sqrt{\frac{\varepsilon}{\mu}} \operatorname{div}\left(\frac{1}{\varepsilon} \operatorname{grad} \mathcal{S}\right)\mathrm{d}s}. \tag{46}$$

Next consider the variation of the complex unit vectors

$$\mathbf{u} = \frac{\mathbf{e}}{\sqrt{\mathbf{e} \cdot \mathbf{e}^\star}}, \qquad \mathbf{v} = \frac{\mathbf{h}}{\sqrt{\mathbf{h} \cdot \mathbf{h}^\star}}, \tag{47}$$

along each ray. Substitution into (41) gives

$$\frac{\partial \mathbf{u}}{\partial \tau} + \frac{1}{2}\left[\frac{\partial \ln (\mathbf{e} \cdot \mathbf{e}^\star)}{\partial \tau} + \nabla^2 \mathcal{S} - \frac{\partial \ln \mu}{\partial \tau}\right]\mathbf{u} + (\mathbf{u} \cdot \operatorname{grad} \ln n)\operatorname{grad} \mathcal{S} = 0.$$

The second, third and fourth terms vanish on account of (43), and it follows that

$$\frac{\mathrm{d}\mathbf{u}}{\mathrm{d}\tau} \equiv n\frac{\mathrm{d}\mathbf{u}}{\mathrm{d}s} = -(\mathbf{u} \cdot \operatorname{grad} \ln n)\operatorname{grad} \mathcal{S}, \tag{48}$$

and similarly

$$\frac{\mathrm{d}\mathbf{v}}{\mathrm{d}\tau} \equiv n\frac{\mathrm{d}\mathbf{v}}{\mathrm{d}s} = -(\mathbf{v} \cdot \operatorname{grad} \ln n)\operatorname{grad} \mathcal{S}. \tag{49}$$

* This relation may also be written in the alternative form

$$\left(\frac{\mathbf{e} \cdot \mathbf{e}^\star}{\mu}\right)_2 = \left(\frac{\mathbf{e} \cdot \mathbf{e}^\star}{\mu}\right)_1 \mathrm{e}^{-\int_{s_1}^{s_2} \frac{\nabla^2 \mathcal{S}}{n} \mathrm{d}s} \tag{45a}$$

which follows when (43) is re-written in the form

$$\frac{\partial}{\partial \tau}\left[\ln\left(\frac{\mathbf{e} \cdot \mathbf{e}^\star}{\mu}\right)\right] = -\nabla^2 \mathcal{S},$$

and the integral is taken along a ray. Eq. (45a) is in fact only another way of expressing the relation (40) for the variation of intensity, and follows from it when the relation

$$I = \frac{2c}{n}\langle w_e \rangle = \frac{c\varepsilon}{8\pi n}(\mathbf{e} \cdot \mathbf{e}^\star)$$

and the Maxwell formula $\varepsilon\mu = n^2$ are used.

This is the required law for the variation of $\mathbf{u}$ and $\mathbf{v}$ along each ray.* In particular, *for a homogeneous medium* ($n =$ constant) (48) and (49) reduce to $\mathrm{d}\mathbf{u}/\mathrm{d}s = \mathrm{d}\mathbf{v}/\mathrm{d}s = 0$ *so that* $\mathbf{u}$ *and* $\mathbf{v}$ *then remain constant along each ray.*

Finally we note that for a time-harmonic homogeneous plane wave in a homogeneous medium, $\mathcal{S} = n\mathbf{s} \cdot \mathbf{r}$ and $\mathbf{e}$, $\mathbf{h}$, $\varepsilon$ and $\mu$ are all constants, and consequently $\mathbf{K} = \mathbf{L} = \mathbf{M} \equiv 0$ in (16). Such a wave (whatever its frequency) therefore obeys rigorously the laws of geometrical optics.

### 3.1.4 Generalizations and the limits of validity of geometrical optics

The considerations of the preceding sections apply to a strictly monochromatic field. Such a field, which may be regarded as a typical Fourier component of an arbitrary field, is produced by a harmonic oscillator, or by a set of such oscillators of the same frequency.

In optics one usually deals with a source which emits light within a narrow, but nevertheless finite, frequency range. The source may then be regarded as arising from a large number of harmonic oscillators whose frequencies fall within this range. To obtain the intensity at a typical field point $P$ one has to sum the individual fields produced by each oscillator (element of the source):

$$\mathbf{E} = \sum_n \mathbf{E}_n, \qquad \mathbf{H} = \sum_n \mathbf{H}_n. \tag{50}$$

The intensity is then given by (using real representation)

$$I(P) = |\langle \mathbf{S} \rangle| = \frac{c}{4\pi} |\langle \mathbf{E} \times \mathbf{H} \rangle| = \frac{c}{4\pi} \left| \sum_{n,m} \langle \mathbf{E}_n \times \mathbf{H}_m \rangle \right|$$

$$= \frac{c}{4\pi} \left| \sum_n \langle \mathbf{E}_n \times \mathbf{H}_n \rangle + \sum_{n \neq m} \langle \mathbf{E}_n \times \mathbf{H}_m \rangle \right|. \tag{51}$$

In many optical problems it is usually permissible to assume that the second sum in (51) vanishes (the fields are then said to be *incoherent*), so that

$$I(P) = \frac{c}{4\pi} \left| \sum_n \langle \mathbf{E}_n \times \mathbf{H}_n \rangle \right| = \left| \sum_n \langle \mathbf{S}_n \rangle \right|, \tag{52}$$

---

* The relations (48) and (49) have an interesting interpretation in terms of non-Euclidean geometry. If we consider the associated non-Euclidean space whose line element is given by

$$\mathrm{d}s' = n\,\mathrm{d}s = n\sqrt{\mathrm{d}x^2 + \mathrm{d}y^2 + \mathrm{d}z^2},$$

then the geometrical light rays correspond to geodesics in this space, and (48) and (49) may be shown to imply that each of the two vectors $\mathbf{u}$ and $\mathbf{v}$ is transferred parallel to itself (in the sense of Levi-Civita parallelism) along each ray. See E. Bortolotti, *Rend. R. Acc. Naz. Linc.*, 6a, **4** (1926), 552; S. M. Rytov, *Compt. Rend. (Doklady) Acad. Sci.* URSS, **18** (1938), 263; R. K. Luneburg, *Mathematical Theory of Optics* (Berkeley and Los Angeles, University of California Press, 1964), pp. 51–55; M. Kline and I. W. Kay, *Electromagnetic Theory and Geometrical Optics* (New York, Interscience Publishers, 1965), pp. 180–183.

More recently the Levi-Civita parallelism has been found to play an important role in connection with the so-called geometrical phase or *Berry phase*, encountered in the transport of some quantum states in parameter space. [See, for example, A Shapere and F. Wilczek (eds.) *Geometric Phases in Physics*, (Singapore, World Scientific, 1989). For a historical account of the subject, see M. Berry, *Physics Today*, Dec. (1990), 34 and I. Bialynicki-Birula, *ibid.*, July (1991), 83.]

$\mathbf{S}_n$ denoting the Poynting vector due to the $n$th element of the source. The more general situation when the second term on the right of (51) has a nonzero value has to be analysed by the use of optical coherence theory (Chapter X).

Let $\delta S$ be a small portion of a wave-front associated with one particular element of the source. Every element of the source sends through $\delta S$ a tube of rays, and the central rays of these tubes fill a cone of solid angle $\delta\Omega$ (Fig. 3.4). If the semivertical angle of this cone is small enough, we may neglect the variation of $\mathbf{S}_n$ with direction, and (52) may then be replaced by

$$I(P) = \sum_n |\langle \mathbf{S}_n \rangle| = \sum_n I_n. \tag{53}$$

Now the number of elements (oscillators) may be regarded as being so large that no appreciable error is introduced by treating the distribution as continuous. The contribution due to each element is then infinitesimal, but the total effect is finite. The sum (integral) is then proportional to $\delta\Omega$:

$$I(P) = B\delta\Omega,$$

and the total (time-averaged) energy flux $\delta F$ which crosses the element $\delta S$ per unit time is given by

$$\delta F = B\delta\Omega\delta S. \tag{54}$$

This formula is of importance in radiometry, and we will return to it in §4.8.1.

We must now briefly consider the limits of validity of geometrical optics. The eikonal equation was derived on the assumption that the terms on the right-hand sides of (11) and (12) may be neglected. If the dimensionless quantities $\varepsilon$, $\mu$ and $|\mathrm{grad}\,\mathcal{S}|$ are assumed to be of order unity, we see that this neglect will be justified provided that the magnitudes of the changes in $\mathbf{e}$ and $\mathbf{h}$ are small compared with the magnitudes of $\mathbf{e}$ and $\mathbf{h}$ over domains whose linear dimensions are of the order of a wavelength. This condition is violated, for example, at boundaries of shadows, for across such boundaries the intensity (and therefore also $\mathbf{e}$ and $\mathbf{h}$) changes rapidly. In the neighbourhood of points where the intensity distribution has a very sharp maximum (e.g. at a focus, see §8.8), geometrical optics likewise cannot be expected to describe correctly the behaviour of the field.

The transport equations (41) and (42) for the complex amplitude vectors $\mathbf{e}$ and $\mathbf{h}$ were obtained on the assumption that $\mathcal{S}$ satisfies the eikonal equation, and that the terms $\lambda_0|\boldsymbol{M}(\mathbf{e}, \varepsilon, \mu)|$ and $\lambda_0|\boldsymbol{M}(\mathbf{h}, \mu, \varepsilon)|$ are small compared with $|\boldsymbol{L}(\mathbf{e}, \mathcal{S}, n, \mu)|$ and $|\boldsymbol{L}(\mathbf{h}, \mathcal{S}, n, \varepsilon)|$, respectively. This imposes certain additional restrictions on, not only
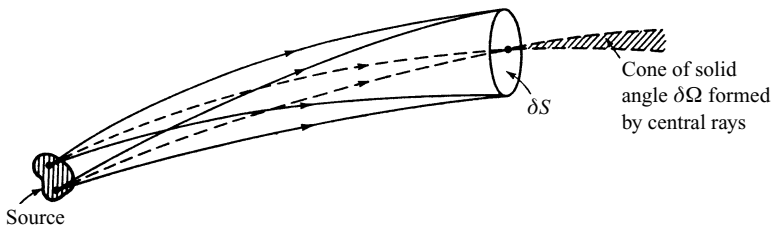


Fig. 3.4 Illustrating the intensity law of geometrical optics with an extended incoherent source.

the first, but also on the second derivatives of **e** and **h**. These conditions are rather complicated and will not be studied here.

It is, of course, possible to obtain improved approximations by retaining some of the higher-order terms in the expansions (18) for the field vectors.[*] In problems of instrumental optics, the practical advantage of such a procedure is, however, doubtful, since the closer the special regions are approached the more terms have to be included, and the expansions usually break down completely at points of particular interest (e.g. at a focus or at a caustic surface). A more powerful approach to the study of the intensity distribution in such regions is offered by methods which will be discussed in the chapters on diffraction.

Finally we stress that the simplicity of the geometrical optics model arises essentially from the fact that, in general, the field behaves *locally* as a plane wave. At optical wavelengths, the regions for which this simple geometrical model is inadequate are an exception rather than a rule; in fact for most optical problems geometrical optics furnishes at least a good starting point for more refined investigations.

## 3.2 General properties of rays

### 3.2.1 The differential equation of light rays

The light rays have been defined as the orthogonal trajectories to the geometrical wave-fronts $\mathcal{S}(x, y, z) = $ constant and we have seen that, if **r** is a position vector of a typical point on a ray and $s$ the length of the ray measured from a fixed point on it, then

$$n \frac{d\mathbf{r}}{ds} = \operatorname{grad} \mathcal{S}. \tag{1}$$

This equation specifies the rays by means of the function $\mathcal{S}$, but one can easily derive from it a differential equation which specifies the rays directly in terms of the refractive index function $n(\mathbf{r})$.

Differentiating (1) with respect to $s$ we obtain

$$\frac{d}{ds}\left(n \frac{d\mathbf{r}}{ds}\right) = \frac{d}{ds}(\operatorname{grad} \mathcal{S})$$

$$= \frac{d\mathbf{r}}{ds} \cdot \operatorname{grad}(\operatorname{grad} \mathcal{S})$$

$$= \frac{1}{n} \operatorname{grad} \mathcal{S} \cdot \operatorname{grad}(\operatorname{grad} \mathcal{S}) \qquad \text{(by (1))}$$

$$= \frac{1}{2n} \operatorname{grad}[(\operatorname{grad} \mathcal{S})^2]$$

$$= \frac{1}{2n} \operatorname{grad} n^2 \qquad \text{[by §3.1 (15)],}$$

---

[*] A theory which takes into account some of the higher-order terms was developed by J. B. Keller and is known as the *geometrical theory of diffraction*. Its central concept is that of a diffracted ray which obeys a generalized Fermat principle. For a review of the theory see J. B. Keller in *J. Opt. Soc. Amer.*, **52** (1962), 116. A collection of some papers on this subject is given in R. C. Hansen (ed.), *Geometrical Theory of Diffraction* (New York, J. Wiley, 1981). Some of its uses are described in G. L. James, *Geometrical Theory of Diffraction for Electromagnetic Waves* (London and New York, The Institute of Electrical Engineers, 1976; revised edition 1980).

i.e.

$$\frac{d}{ds}\left(n\frac{d\mathbf{r}}{ds}\right) = \text{grad } n. \tag{2}$$

This is the vector form of the differential equations of the light rays. In particular, *in a homogeneous medium n* = constant and (2) then reduces to

$$\frac{d^2\mathbf{r}}{ds^2} = 0,$$

whence

$$\mathbf{r} = s\mathbf{a} + \mathbf{b}, \tag{3}$$

**a** and **b** being constant vectors. Eq. (3) is a vector equation of a straight line in the direction of the vector **a**, passing through the point $\mathbf{r} = \mathbf{b}$. Hence *in a homogeneous medium the light rays have the form of straight lines*.

As an example of some interest, let us consider rays in a medium which has spherical symmetry, i.e. where the refractive index depends only on the distance $r$ from a fixed point $O$:

$$n = n(r). \tag{4}$$

This case is approximately realized by the earth's atmosphere, when the curvature of the earth is taken into account.

Consider the variation of the vector $\mathbf{r} \times [n(\mathbf{r})\mathbf{s}]$ along the ray. We have

$$\frac{d}{ds}(\mathbf{r} \times n\mathbf{s}) = \frac{d\mathbf{r}}{ds} \times n\mathbf{s} + \mathbf{r} \times \frac{d}{ds}(n\mathbf{s}). \tag{5}$$

Since $d\mathbf{r}/ds = \mathbf{s}$, the first term on the right vanishes. The second term may, on account of (2), be written as $\mathbf{r} \times \text{grad } n$. Now from (4)

$$\text{grad } n = \frac{\mathbf{r}}{r}\frac{dn}{dr},$$

so that the second term on the right-hand side of (5) also vanishes. Hence

$$\mathbf{r} \times n\mathbf{s} = \text{constant}. \tag{6}$$

This relation implies that all the rays are plane curves, situated in a plane through the origin, and that along each ray

$$nr\sin\phi = \text{constant}, \tag{7}$$

where $\phi$ is the angle between the position vector **r** and the tangent at the point **r** on the ray (see Fig. 3.5). Since $r\sin\phi$ represents the perpendicular distance $d$ from the origin to the tangent, (7) may also be written as

$$nd = \text{constant}. \tag{8}$$

This relation is sometimes called the *formula of Bouguer* and is the analogue of a well-known formula in dynamics, which expresses the conservation of angular momentum of a particle moving under the action of a central force.

To obtain an explicit expression for the rays in a spherically symmetrical medium, we recall from elementary geometry that, if $(r, \theta)$ are the polar coordinates of a plane
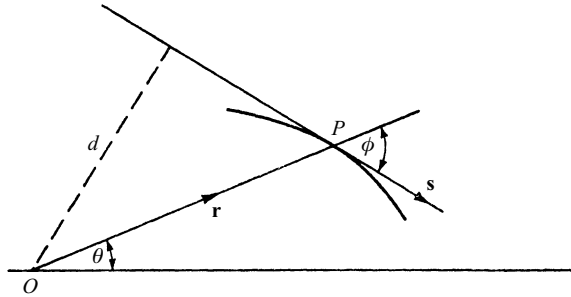
Fig. 3.5 Illustrating Bouguer's formula, $nd = $ constant, for rays in a medium with spherical symmetry.

curve, then the angle $\phi$ between the radius vector to a point $P$ on the curve and the tangent at $P$ is given by[*]

$$\sin \phi = \frac{r(\theta)}{\sqrt{r^2(\theta) + \left(\dfrac{\mathrm{d}r}{\mathrm{d}\theta}\right)^2}}.$$  (9)

From (7) and (9)

$$\frac{\mathrm{d}r}{\mathrm{d}\theta} = \frac{r}{c}\sqrt{n^2 r^2 - c^2},$$  (10)

$c$ being a constant. The equation of rays in a medium with spherical symmetry may therefore be written in the form

$$\theta = c \int^r \frac{\mathrm{d}r}{r\sqrt{n^2 r^2 - c^2}}.$$  (11)

Let us now return to the general case and consider the *curvature vector* of a ray, i.e. the vector

$$\mathbf{K} = \frac{\mathrm{d}\mathbf{s}}{\mathrm{d}s} = \frac{1}{\rho}\boldsymbol{v},$$  (12)

whose magnitude $1/\rho$ is the reciprocal of the radius of curvature; $\boldsymbol{v}$ is the unit principal normal at a typical point of the ray.

From (2) and (12) it follows that

$$n\mathbf{K} = \operatorname{grad} n - \frac{\mathrm{d}n}{\mathrm{d}s}\mathbf{s}.$$  (13)

This relation shows that *the gradient of the refractive index lies in the osculating plane of the ray.*

If we multiply (13) scalarly by $\mathbf{K}$ and use (12) we find that

$$|\mathbf{K}| = \frac{1}{\rho} = \boldsymbol{v} \cdot \operatorname{grad} \ln n.$$  (14)

[*] See, for example, R. Courant, *Differential and Integral Calculus*, Vol. I (Glasgow, Blackie, 2nd edition, 1942), p. 265.
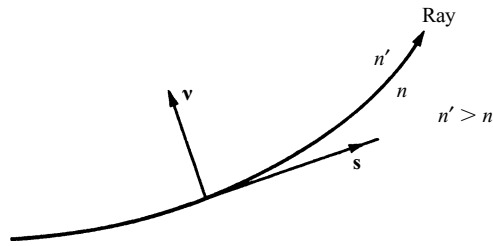
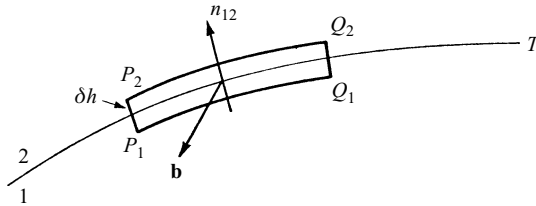Fig. 3.6 Bending a ray in a heterogeneous medium.



Fig. 3.7 Illustration of the laws of refraction and reflection.

Since $\rho$ is always positive, this implies that as we proceed along the principal normal the refractive index increases, i.e. *the ray bends towards the region of higher refractive index* (Fig. 3.6).

### 3.2.2  The laws of refraction and reflection

So far it has been assumed that the refractive index function $n$ is continuous. We must now discuss the behaviour of rays when they cross a surface separating two homogeneous media of different refractive indices. It has been shown by Sommerfeld and Runge[*] that the behaviour can easily be determined by an argument similar to that used in deriving the conditions relating to the changes in the field vectors across a surface discontinuity (see §1.1.3).

It follows from (1), on account of the identity $\operatorname{curl} \operatorname{grad} \equiv 0$, that the vector $n\mathbf{s} = n\,\mathrm{d}\mathbf{r}/\mathrm{d}s$, called sometimes the *ray vector*, satisfies the relation

$$\operatorname{curl} n\mathbf{s} = 0. \tag{15}$$

As in §1.1.3 we replace the discontinuity surface $T$ by a transition layer throughout which $\varepsilon$, $\mu$ and $n$ change rapidly but continuously from their values near $T$ on one side to their values near $T$ on the other. Next we take a plane element of area with its sides $P_1 Q_1$ and $P_2 Q_2$ parallel and with $P_1 P_2$ and $Q_1 Q_2$ perpendicular to $T$ (Fig. 3.7). If $\mathbf{b}$ denotes the unit normal to this area, then we have from (15), on integrating throughout the area and applying Stokes' theorem,

$$\int (\operatorname{curl} n\mathbf{s}) \cdot \mathbf{b}\,\mathrm{d}S = \int n\mathbf{s} \cdot \mathrm{d}\mathbf{r} = 0, \tag{16}$$

[*]  Sommerfeld and Runge, *loc cit.*

the second integral being taken along the boundary curve $P_1 Q_1 Q_2 P_2$. Proceeding to the limit as the height $\delta h \to 0$, in a strictly similar manner as in the derivation of §1.1 (23), we obtain

$$\mathbf{n}_{12} \times (n_2 \mathbf{s}_2 - n_1 \mathbf{s}_1) = 0, \tag{17}$$

where $\mathbf{n}_{12}$ is the unit normal to the boundary surface pointing from the first into the second medium. Eq. (17) implies that *the tangential component of the ray vector $n\mathbf{s}$ is continuous across the surface* or, what amounts to the same thing, *the vector $\mathbf{N}_{12} = n_2 \mathbf{s}_2 - n_1 \mathbf{s}_1$ is normal to the surface.*

Let $\theta_1$ and $\theta_2$ be the angles which the incident ray and the refracted ray make with the normal $\mathbf{n}_{12}$ to the surface (see Fig. 3.8(a)). Then it follows from (17) that

$$n_2(\mathbf{n}_{12} \times \mathbf{s}_2) = n_1(\mathbf{n}_{12} \times \mathbf{s}_1), \tag{18}$$

so that

$$n_2 \sin \theta_2 = n_1 \sin \theta_1. \tag{19}$$

Eq. (18) implies that *the refracted ray lies in the same plane as the incident ray and the normal to the surface* (*the plane of incidence*) and (19) shows that *the ratio of the sine of the angle of refraction to the sine of the angle of incidence is equal to the ratio $n_1/n_2$ of the refractive indices.* These two results express *the law of refraction* (*Snell's law*). This law has already been derived in §1.5 for the special case of plane waves. But whilst the earlier discussion concerned a plane wave of *arbitrary* wavelength falling upon a plane refracting surface, the present analysis applies to waves and refracting surfaces of more general form, provided that the wavelength is sufficiently small ($\lambda_0 \to 0$). This condition means, in practice, that the radii of curvature of the incident wave and of the boundary surface must be large compared to the wavelength of the incident light.

As in the case treated in §1.5 we must expect that there will be another wave, the reflected wave, propagated back into the first medium. Setting $n_2 = n_1$ in (18) and (19) [see Fig. 3.8(b)] it follows that *the reflected ray lies in the plane of incidence* and that $\sin \theta_2 = \sin \theta_1$; hence

$$\theta_2 = \pi - \theta_1. \tag{20}$$

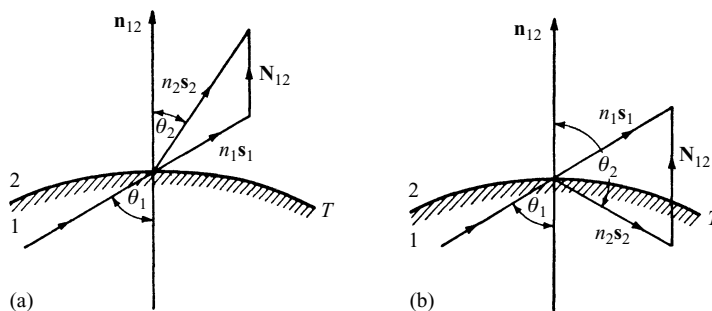The last two results express *the law of reflection*.



Fig. 3.8 Illustrating the laws of (a) refraction and (b) reflection.

### 3.2.3 Ray congruences and their focal properties

The relation (15), namely

$$\operatorname{curl} n\mathbf{s} = 0, \tag{21}$$

characterizes all the ray systems which can be realized in an isotropic medium and distinguishes them from more general families of curves. In a homogeneous isotropic medium $n$ is constant, and (21) then reduces to

$$\operatorname{curl} \mathbf{s} = 0. \tag{22}$$

Rays in a heterogeneous isotropic medium can also be characterized by a relation independent of $n$. It may be obtained by applying to (21) the identity $\operatorname{curl} n\mathbf{s} = n \operatorname{curl} \mathbf{s} + (\operatorname{grad} n) \times \mathbf{s}$ and taking the scalar product with $\mathbf{s}$. It then follows that a system of rays in any *isotropic* medium must satisfy the relation

$$\mathbf{s} \cdot \operatorname{curl} \mathbf{s} = 0. \tag{23}$$

A system of curves which fills a portion of space in such a way that in general a single curve passes through each point of the region is called a *congruence*. If there exists a family of surfaces which cut each of the curves orthogonally the congruence is said to be *normal*; if there is no such family, it is said to be *skew*. For ordinary geometrical optics (light propagation) only normal congruences are of interest, but in electron optics (see Appendix II) skew congruences also play an important part.

If each curve of the congruence is a straight line the congruence is said to be *rectilinear*; (23) and (22) are the necessary and sufficient conditions that the curves should represent a *normal* and a *normal rectilinear congruence*, respectively.[*]

Let us choose a set of curvilinear coordinate lines $u$, $v$ on one of the orthogonal surfaces $\mathcal{S}(x, y, z) = \text{constant}$. To every point $Q(u, v)$ of this surface there will then correspond one curve of the congruence, namely that curve which meets $\mathcal{S}$ in $Q$. Let $\mathbf{r}$ denote the position vector of a point $P$ on the curve. $\mathbf{r}$ may then be regarded as a function of the coordinates $(u, v)$ and of the length of arc $s$ between $Q$ and $P$, measured along the curve (Fig. 3.9).

Consider two neighbouring curves of the congruence passing through the points
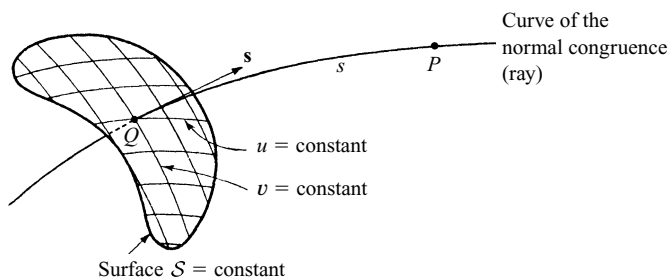


Fig. 3.9 Notation relating to normal congruence.

---

[*] For a more detailed discussion of congruences of curves see, for example, C. E. Weatherburn, *Differential Geometry of Three Dimensions* (Cambridge, Cambridge University Press), Vol. I (1927), Chapter X; Vol. II (1930), Chapter XIII.

$(u, v)$ and $(u + \mathrm{d}u, v + \mathrm{d}v)$ on $\mathcal{S}$, and let us examine whether there are points on these curves such that the distance between them is of the second or higher order (one says that the curves cut to first order at such points). Points with this property are called *foci* and must satisfy the equation

$$\mathbf{r}(u, v, s) = \mathbf{r}(u + \mathrm{d}u, v + \mathrm{d}v, s + \mathrm{d}s) \tag{24}$$

to the first order.

Expanding (24) we obtain

$$\mathbf{r}_u \, \mathrm{d}u + \mathbf{r}_v \, \mathrm{d}v + \mathbf{s} \, \mathrm{d}s = 0, \tag{25}$$

where $\mathbf{r}_u$, $\mathbf{r}_v$ are the partial derivatives with respect to $u$ and $v$. Condition (25) implies that $\mathbf{r}_u$, $\mathbf{r}_v$ and $\mathbf{s}$ are coplanar. This is equivalent to saying that the scalar triple product of the three vectors vanishes, i.e.

$$[\mathbf{r}_u, \mathbf{r}_v, \mathbf{s}] = 0. \tag{26}$$

The number of foci on a given curve $(u, v)$ depends on the number of values of $s$ which satisfy (26). If $\mathbf{r}$ is a polynomial in $s$ of degree $m$, then since $\mathbf{s} = \mathrm{d}\mathbf{r}/\mathrm{d}s$, it is seen that (26) is an equation of degree $3m - 1$ in $s$. In particular, if the congruence is rectilinear, $\mathbf{r}$ is a linear function of $s$ ($m = 1$), showing that *there are two foci on each ray of a rectilinear congruence.*

If $u$ and $v$ take on all possible values, the foci will describe a surface, represented by (26), known as the *focal surface*; in optics it is called the *caustic surface*. Any curve of the congruence is tangent to the focal surface at each focus of the curve. The tangent plane at any point of the focal surface is known as the *focal plane*.

We shall mainly be concerned with rays in a homogeneous medium, i.e. with rectilinear congruences. Some further properties of such congruences will be discussed in §4.6, in connection with astigmatic pencils of rays.

## 3.3 Other basic theorems of geometrical optics

With the help of the relations established in the preceding sections, we shall now derive a number of theorems concerning rays and wave-fronts.

### 3.3.1 Lagrange's integral invariant

Assume first that the refractive index $n$ is a continuous function of position. Then as in §3.2 (16) it follows on applying Stokes' theorem to the integral, taken over any open surface, of the normal component of curl $n\mathbf{s}$, that

$$\oint n\mathbf{s} \cdot \mathrm{d}\mathbf{r} = 0. \tag{1}$$

The integral extends over the closed boundary curve $C$ of the surface. Eq. (1) is known as *Lagrange's integral invariant*[*] and implies that *the integral*

---

[*] Sometimes called *Poincaré's invariant*. In fact it is only a special one-dimensional case of much more general integral invariants discussed by J. H. Poincaré in his *Les Méthodes Nouvelles de la Mécanique Céleste*, Vol. 3 (Paris, Gauthier-Villars, 1899). See E. Cartan, *Leçons sur les Invariants Integraux* (Paris, Hermann, 1922). See also our Appendix I, eq. (85).
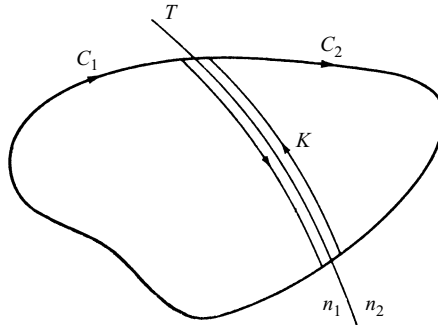
Fig. 3.10 Illustrating the derivation of Lagrange's integral invariant in the presence of a surface discontinuity of the refractive index.

$$\int_{P_1}^{P_2} n\mathbf{s} \cdot \mathrm{d}\mathbf{r} \tag{2}$$

*taken between any two points $P_1$ and $P_2$ in the field, is independent of the path of integration.*

   With the help of the law of refraction it is easily seen that (1) also holds when the curve $C$ intersects a surface which separates two homogeneous media of different refractive indices. To show this, let $C_1$ and $C_2$ be the portions of $C$ on each side of the refracting surface $T$ (Fig. 3.10), and let the points of intersection of $C$ with the surface $T$ be joined by another curve $K$ in the surface. On taking (1) along each of the loops $C_1K$ and $C_2K$ and on adding the equations, we obtain

$$\int_{C_1} n_1\mathbf{s}_1 \cdot \mathrm{d}\mathbf{r} + \int_{C_2} n_2\mathbf{s}_2 \cdot \mathrm{d}\mathbf{r} + \int_K (n_2\mathbf{s}_2 - n_1\mathbf{s}_1) \cdot \mathrm{d}\mathbf{r} = 0. \tag{3}$$

The integral over $K$ vanishes, since according to the law of refraction the vector $\mathbf{N}_{12} = n_1\mathbf{s}_1 - n_2\mathbf{s}_2$ is at each point of $K$ perpendicular to the surface, and consequently (3) reduces to (1).

### 3.3.2  The principle of Fermat

The *principle of Fermat*, known also as the principle of the *shortest optical path*,[*] asserts that *the optical length*

$$\int_{P_1}^{P_2} n \, \mathrm{d}s \tag{4}$$

*of an actual ray between any two points $P_1$ and $P_2$ is shorter than the optical length of any other curve which joins these points and which lies in a certain regular neighbourhood of it.* By a regular neighbourhood we mean one that may be covered by

---

[*] Since by §3.1 (27)

$$\int_{P_1}^{P_2} n \, \mathrm{d}s = c \int_{P_1}^{P_2} \mathrm{d}t$$

it is also known as the *principle of least time.*

rays in such a way that one (and only one) ray passes through each point of it. Such a covering is exhibited, for example, by rays from a point source $P_1$ in that domain around $P_1$ where the rays on account of refraction or reflection or on account of their curvature do not intersect each other.

Before proving this theorem it may be mentioned that it is possible to formulate Fermat's principle in a form which is weaker but which has a wider range of validity. According to this formulation the actual ray is distinguished from other curves (no longer restricted to lie in a regular neighbourhood) by a *stationary value* of the integral.*

To prove Fermat's principle, we take a pencil of rays and compare a segment $P_1 P_2$ of a ray $\overline{C}$ with an arbitrary curve $C$ joining $P_1$ and $P_2$ (Fig. 3.11). Let two neighbouring orthogonal trajectories (wave-fronts) of the pencil intersect $C$ in $Q_1$ and $Q_2$ and $\overline{C}$ in $\overline{Q}_1$ and $\overline{Q}_2$. Further let $Q_2'$ be the point of intersection of the trajectory $Q_2\overline{Q}_2$ with the ray $\overline{C}'$ which passes through $Q_1$.

Applying Lagrange's integral relation to the small triangle $Q_1 Q_2 Q_2'$, we have

$$(n\mathbf{s} \cdot \mathrm{d}\mathbf{r})_{Q_1 Q_2} + (n\mathbf{s} \cdot \mathrm{d}\mathbf{r})_{Q_2 Q_2'} - (n\,\mathrm{d}s)_{Q_1 Q_2} = 0. \tag{5}$$

Now from the definition of the scalar product

$$(n\mathbf{s} \cdot \mathrm{d}\mathbf{r})_{Q_1 Q_2} \leqslant (n\,\mathrm{d}s)_{Q_1 Q_2}.$$

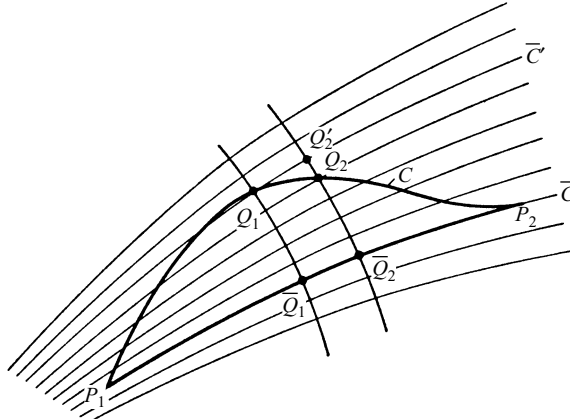Further, $\mathbf{s}$ is orthogonal to $\mathrm{d}\mathbf{r}$ on the wave-front, so that



Fig. 3.11 Illustrating Fermat's principle.

---

* To find the curves for which the integral has a stationary value we must apply in general the methods of the variational calculus, described in Appendix I. It is shown there that such curves satisfy the Euler differential equations Appendix I (7). In the present case these are nothing but the equations §3.2 (2) of the rays as shown in §11 of Appendix I.

It has been stressed by C. Carathéodory [*Geometrische Optik* (Berlin, Springer, 1937)] that the stationary value is never a true maximum. In the weaker formulation of Fermat's principle it is therefore appropriate to speak of a stationary value rather than of an extremal value. The minimal formulation on the other hand corresponds to a 'strong minimum' in the sense of Jacobi (Appendix I, §10).

$$(n\mathbf{s} \cdot \mathbf{dr})_{Q_2 Q_2'} = 0.$$

Also from §3.1 (25), since $Q_1$, $Q_2'$ and $\overline{Q}_1$, $\overline{Q}_2$ are corresponding points on the two wave-fronts,

$$(n\,ds)_{Q_1 Q_2'} = (n\,ds)_{\overline{Q}_1 \overline{Q}_2}.$$

On substituting from the last three relations into (5) we find that

$$(n\,ds)_{\overline{Q}_1 \overline{Q}_2} \lessgtr (n\,ds)_{Q_1 Q_2}, \tag{6}$$

whence, on integration,

$$\int_{\overline{C}} n\,ds \lessgtr \int_{C} n\,ds. \tag{7}$$

Moreover, the equality sign could only hold if the directions of $\mathbf{s}$ and $\mathbf{dr}$ were coincident at every point of $C$, i.e. if the comparison curve was an actual ray. This case is excluded by our assumption that not more than one ray passes through any point of the neighbourhood. Hence the optical length of the ray is smaller than the optical length of the comparison curve, which is Fermat's principle.

It can easily be seen that, when the regularity condition is not fulfilled, the optical length of the ray may no longer be a minimum. Consider for example a field of rays from a point source $P_1$ in a homogeneous medium, reflected by a plane mirror (Fig. 3.12). Two rays then pass through each point $P_2$; the optical length of the direct ray $P_1 P_2$ is an absolute minimum but the reflected ray $P_1 M P_2$ gives a minimum only relative to curves in a certain restricted neighbourhood of it. In general when rays from a point source $P_1$ are refracted or reflected at boundaries between homogeneous media, the regular neighbourhood will terminate on the envelope (caustic) formed by the rays. The point $P_1'$ at which a ray from a point source at $P_1$ touches the envelope is called the *conjugate* of $P_1$ on the particular ray. For the optical length of a ray $P_1 P_2$ to be a minimum, $P_2$ must lie between $P_1$ and $P_1'$, i.e. $P_1$ and $P_2$ must lie on the same side of the caustic. For example, in the case of an uncorrected lens (Fig. 3.13) the central ray from $P_1$ has a minimal optical length only up to the tip ($P_1'$) of the caustic (the Gaussian image of $P_1$). For any point $P_2$ which lies behind the envelope the optical length of the direct path $P_1 P_1' P_2$ exceeds that of the broken path $P_1 A B P_2$.
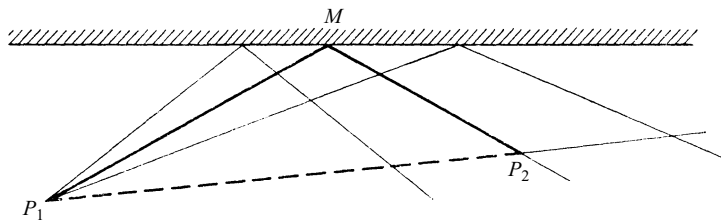


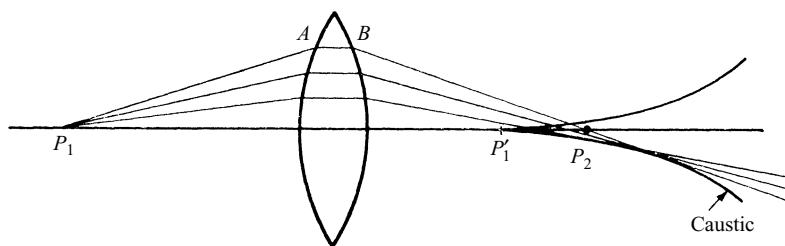Fig. 3.12 Field of rays obtained by a reflection of light from a point source on a plane mirror.

Fig. 3.13 Caustic formed by rays from an axial point source after passing through a lens.

### 3.3.3 The theorem of Malus and Dupin and some related theorems

The light rays have been defined as the orthogonal trajectories of the wave surfaces $\mathcal{S}(x, y, z) = $ constant, $\mathcal{S}$ being a solution of the eikonal equation §3.1 (15). This is a natural way of introducing the light rays when the laws of geometrical optics are to be deduced from Maxwell's equations. Historically, however, geometrical optics developed as the theory of light rays which were defined differently, namely as curves for which the line integral $\int n\,\mathrm{d}s$ has a stationary value. Formulated this way geometrical optics may then be developed purely along the lines of calculus of variations.[*]

Variational considerations are of considerable importance as they often reveal analogies between different branches of physics. In particular there is a close analogy between geometrical optics and the mechanics of a moving particle; this was brought out very clearly by the celebrated investigations of Sir W. R. Hamilton, whose approach became of great value in modern physics, especially in applications to De Broglie's wave mechanics. In order not to interrupt the optical considerations, accounts of the relevant parts of the calculus of variations and of the Hamiltonian analogy are given in separate sections (Appendices I and II). Here we shall only show how several theorems, which played an important part in the development of geometrical optics, may be derived from Lagrange's integral invariant.

Consider rays in a homogeneous medium: if they all have a point in common, e.g. when they then proceed from a point source, they are said to form a *homocentric pencil*. Such a pencil forms a normal congruence, since every ray of the pencil is cut orthogonally by spheres centred on the mutual point of intersection of the rays. In 1808 Malus[†] showed that, if a homocentric pencil of rectilinear rays is refracted or reflected at a surface, the resulting pencil (in general no longer homocentric) will again form a normal congruence. Later Dupin (1816), Quetelet (1825), and Gergonne (1825) generalized Malus's result. These investigations lead to the following theorem, known sometimes as *the theorem of Malus and Dupin: A normal rectilinear congruence remains normal after any number of refractions or reflections*.[‡]

---

[*] A systematic treatment of this kind is given for example in C. Carathéodory (*loc. cit.*).

[†] E. Malus, Optique Dioptrique, *J. École polytechn.*, **7** (1808), 1–44, 84–129. Also his Traité d'optique, *Mém. présent. à l'Institut par divers savants*, **2** (1811), 214–302. References and an account of the interesting history of the Malus–Dupin theorem can be found in the *Mathematical Papers of Sir William Rowan Hamilton*, Vol. 1 (*Geometrical Optics*), eds. A. W. Conway and J. L. Synge (Cambridge, Cambridge University Press, 1931), p. 463.

[‡] T. Levi-Civita, *Rend. R. Acc. Naz. Linc.*, **9** (1900) 237 established the converse theorem, namely that in general two normal rectilinear congruences may be transformed into each other by a single refraction or reflection.

It will be sufficient to establish the theorem for a single refraction. Consider a normal rectilinear congruence of rays in a homogeneous medium of refractive index $n_1$ and assume that the rays undergo a refraction at a surface $T$ which separates this medium from another homogeneous medium of refractive index $n_2$ (Fig. 3.14).

Let $S_1$ be one of the orthogonal trajectories (wave-fronts) in the first region, and let $A_1$ and $P$ be the points of intersections of a typical ray in the first medium with $S_1$ and with $T$ respectively, and let $A_2$ be any point on the refracted ray. If the point $A_1$ is displaced to another point $B_1$ on the wave-front, the point $P$ will be displaced to another point $Q$ on the refracting surface. Now take a point $B_2$, on the ray which is refracted at $Q$, such that the optical path from $B_1$ to $B_2$ is equal to the optical path from $A_1$ to $A_2$:

$$[A_1 P A_2] = [B_1 Q B_2]. \tag{8}$$

As $B_1$ takes on all possible positions on $S_1$ the point $B_2$ describes a surface $S_2$. It will now be shown that the refracted ray $Q B_2$ is perpendicular to this surface.

Applying Lagrange's integral invariant to the closed path $A_1 P A_2 B_2 Q B_1 A_1$, it follows that

$$\int_{A_1 P A_2} n\,\mathrm{d}s + \int_{A_2 B_2} n\mathbf{s} \cdot \mathrm{d}\mathbf{r} + \int_{B_2 Q B_1} n\,\mathrm{d}s + \int_{B_1 A_1} n\mathbf{s} \cdot \mathrm{d}\mathbf{r} = 0. \tag{9}$$

Now by (8),

$$\int_{A_1 P A_2} n\,\mathrm{d}s + \int_{B_2 Q B_1} n\,\mathrm{d}s = 0. \tag{10}$$

Moreover, since on $S_1$ the unit vector $\mathbf{s}$ is everywhere orthogonal to $S_1$,

$$\int_{B_1 A_1} n\mathbf{s} \cdot \mathrm{d}\mathbf{r} = 0, \tag{11}$$

so that (9) reduces to

$$\int_{A_2 B_2} n\mathbf{s} \cdot \mathrm{d}\mathbf{r} = 0. \tag{12}$$

This relation must hold for every curve on $S_2$. This is only possible if $\mathbf{s} \cdot \mathrm{d}\mathbf{r} = 0$ for every linear element $\mathrm{d}\mathbf{r}$ of $S_2$, i.e. if the refracted rays are orthogonal to the surface; in
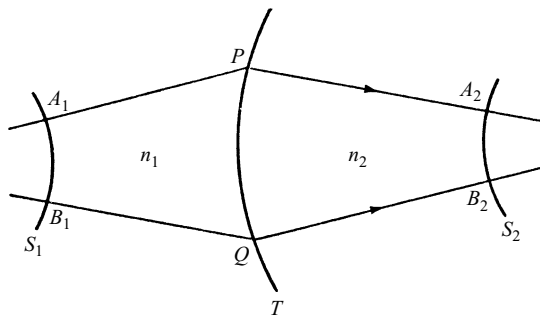


Fig. 3.14 Illustrating the theorem of Malus and Dupin.

other words *if the refracted rays form a normal congruence*. The proof for reflection is strictly analogous.

Since $[A_1PA_2] = [B_1QB_2]$ it follows that *the optical path length between any two orthogonal surfaces (wave-fronts) is the same for all rays*. This result clearly remains valid when several successive refractions or reflections take place and, as is immediately obvious from §3.1 (26) it also applies to rays in a medium with a continuously varying refractive index. This theorem is known as the *principle of equal optical path*; it implies that the orthogonal trajectories (geometrical wave-fronts) of a normal congruence of rays, or of a set of normal congruences generated by successive refractions or reflections, are 'optically parallel' to each other (see Appendix I).

A related theorem, first put forward by Huygens[*] asserts that *each element of a wave-front may be regarded as the centre of a secondary disturbance which gives rise to spherical wavelets*; and, moreover, that *the position of the wave-front at any later time is the envelope of all such wavelets*. This result, sometimes called *Huygens' construction*, is essentially a rule for the construction of a set of surfaces which are 'optically parallel' to each other. If the medium is homogeneous, one can use in the construction wavelets of finite radius, in other cases one has to proceed in infinitesimal steps.

Huygens' theorem was later extended by Fresnel and led to the formulation of the so-called Huygens–Fresnel principle, which is of great importance in the theory of diffraction (see §8.2), and which may be regarded as the basic postulate of the wave theory of light.

---

[*] Chr. Huygens, *Traité de la Lumière* (Leyden, 1690); English translation (*Treatise on Light*) by S. P. Thompson (London, Macmillan & Co., 1912).