**Chapter 1** 

# Basic principles of the ionosphere

# 1.1 Introduction

# 1.1.1 The ionosphere and radio-wave propagation

The *ionosphere* is the ionized component of the atmosphere, comprising free electrons and positive ions, generally in equal numbers, in a medium that is electrically neutral. Though the charged particles are only a minority amongst the neutral ones, they nevertheless exert a great influence on the electrical properties of the medium, and it is their presence that brings about the possibility of radio communication over large distances by making use of one or more ionospheric reflections.

The early history of the ionosphere is very much bound up with the development of communications. The first suggestions that there are electrified layers within the upper atmosphere go back to the nineteenth century, but the modern developments really started with Marconi's well-known experiments in trans-Atlantic communication (from Cornwall to Newfoundland) in 1901. These led to the suggestions by Kennelly and by Heaviside (made independently) that, because of the Earth's curvature, the waves could not have traveled directly across the Atlantic but must have been reflected from an ionized layer. The name *ionosphere* came into use about 1932, having been coined by Watson-Watt several years previously. Subsequent research has revealed a great deal of information about the ionosphere: its vertical structure, its temporal and spatial variations, and the physical processes by which it is formed and which influence its behavior.

Looked at most simply, the ionosphere acts as a mirror situated between 100 and 400 km above the Earth's surface, as in Figure 1.1, which allows reflected

1

Basic principles of the ionosphere



Figure 1.1. Long distance propagation by multiple hops between the ionosphere and the ground.

signals to reach points around the bulge of the Earth. The details of how reflection occurs depend on the radio frequency of the signal, but the most usual mechanism, which applies in the high-frequency (HF) band (3–30 MHz), is actually a gradual bending of the ray towards the horizontal as the refractive index of the ionospheric medium decreases with altitude. Under good conditions, signals can be propagated in this way for several thousand kilometers by means of repeated reflections between ionosphere and ground. Reflection from a higher level (the F region) obviously gives a greater range per "hop" than does one from a lower level (the E region), but which mode is possible depends on the structure of the ionosphere at the time. Higher radio frequencies tend to be reflected from greater heights, but if the frequency is too high there may be insufficient bending and the signal then penetrates the layer and is lost to space. This is the first complication of radio propagation.

The second complication is that the lower layers of the ionosphere tend to absorb the signal. This effect is greater for signals of lower frequency and greater obliquity. Hence, practical radio communication generally requires a compromise. The ionosphere is constantly changing, and the art of propagation prediction is to determine the best radio frequency for a given path for the current state of the ionosphere. Plainly, an understanding of ionospheric mechanisms is basic to efficient radio communication.

Further details about radio propagation are given in Chapter 3, and our central topic of how propagation at high latitudes is affected by the vagaries of the high-latitude ionosphere is discussed later in the book.

## 1.1.2 Why the ionosphere is so different at high latitude

The terrestrial ionosphere may be divided broadly into three regions that have rather different properties according to their geomagnetic latitude. The midlatitude region has been explored the most completely and is the best understood. There, the ionization is produced almost entirely by energetic ultra-violet and Xray emissions from the Sun, and is removed again by chemical recombination processes that may involve the neutral atmosphere as well as the ionized species. The

#### 1.1 Introduction

movement of ions, and the balance between production and loss, are affected by winds in the neutral air. The processes typical of the mid-latitude ionosphere also operate at high and low latitudes, but in those regions additional processes are also important.

The low-latitude zone, spanning 20° or 30° either side of the magnetic equator, is strongly influenced by electromagnetic forces that arise because the geomagnetic field runs horizontally over the magnetic equator. The primary consequence is that the electrical conductivity is abnormally large over the equator. A strong electric current (an "electrojet") flows in the E region, and the F region is subject to electrodynamic lifting and a "fountain effect" that distorts the general form of the ionosphere throughout the low-latitude zone.

At high latitudes we find the opposite situation. Here the geomagnetic field runs nearly vertical, and this simple fact of nature leads to the existence of an ionosphere that is considerably more complex than that in either the middle or the low-latitude zones. This happens because the magnetic field-lines connect the high latitudes to the outer part of the magnetosphere which is driven by the solar wind, whereas the ionosphere at middle latitude is connected to the inner magnetosphere, which essentially rotates with the Earth and so is less sensitive to external influence. We can immediately identify four general consequences.

- (a). The high-latitude ionosphere is dynamic. It circulates in a pattern mainly controlled by the solar wind but which is also variable.
- (b). The region is generally more accessible to energetic particle emissions from the Sun that produce additional ionization. Thus it is affected by sporadic events, which can seriously degrade polar radio propagation. Over a limited range of latitudes the dayside ionosphere is directly accessible to material from the solar wind.
- (c). The auroral zones occur within the high-latitude region. Again, their location depends on the linkage with the magnetosphere, in this case into the distorted tail of the magnetosphere. The auroral phenomena include electrojets, which cause magnetic perturbations, and there are "substorms" in which the rate of ionization is greatly increased by the arrival of energetic electrons. The auroral regions are particularly complex for radio propagation.
- (d). A "trough" of lesser ionization may be formed between the auroral and the mid-latitude ionospheres. Although the mechanisms leading to the formation of the trough are not completely known, it is clear that one fundamental cause is the difference in circulation pattern between the inner and outer parts of the magnetosphere.

This monograph is concerned mainly with the ionosphere at high latitudes, but before considering the special behavior which occurs in those regions we must review some processes affecting the ionosphere in general and summarize the more normal behavior at middle latitudes. In order to do that, we must first

Basic principles of the ionosphere



**Figure 1.2.** Nomenclature of the upper atmosphere based on temperature, composition, mixing, and ionization. (J. K. Hargreaves, *The Solar–Terrestrial Environment*. Cambridge University Press, 1992.)

consider the nature of the neutral upper atmosphere in which the ionosphere is formed.

### 1.2 The vertical structure of the atmosphere

### 1.2.1 Nomenclature

A static planetary atmosphere may be described by four properties: pressure (P), density ( $\rho$ ), temperature (T), and composition. Since these are not independent it is not necessary to specify all of them. The nomenclature of the atmosphere is based principally on the variation of temperature with height, as in Figure 1.2. Here, the different regions are called "spheres" and the boundaries between them are "pauses". The lowest region is the troposphere, in which the temperature falls off with increasing height at a rate of 10 K km<sup>-1</sup> or less. Its upper boundary is the tropopause at a height of 10-12 km. The stratosphere which lies above it was once thought to be isothermal, but it is actually a region where the temperature increases with height. At about 50 km is a maximum due to the absorption of solar ultra-violet radiation in ozone; this is the stratopause. Above that the temperature again decreases in the mesosphere (or middle atmosphere) and passes through another minimum at the mesopause at 80-85 km. At about 180 K, this is the coldest part of the whole atmosphere. Above the mesopause, heating by solar ultra-violet radiation ensures that the temperature gradient remains positive, and this is the thermosphere. Eventually the temperature of the thermosphere becomes

#### 1.2 Vertical structure

almost constant at a value that varies with time but is generally over 1000 K; this is the hottest part of the atmosphere.

Though the classification by temperature is generally the most useful, others based on the state of mixing, the composition or the state of ionization are also useful. The lowest part of the atmosphere is well mixed, with a composition much like that at sea level except for minor components. This is the *turbosphere* or *homosphere*. In the upper region, essentially the thermosphere, mixing is inhibited by the positive temperature gradient, and here, in the *heterosphere*, the various components separate under gravity and as a result the composition varies with altitude. The boundary between the two regions, which occurs at about 100 km, is the *turbopause*. Above the turbopause the gases separate by gaseous diffusion more rapidly than they are mixed by turbulence.

Within the heterosphere there are regions where helium or hydrogen may be the main component. These are the *heliosphere* and the *protonosphere*, respectively. From the higher levels, above about 600 km, individual atoms can escape from the Earth's gravitational attraction; this region is called the *exosphere*. The base of the exosphere is the *exobase* or the *baropause*, and the region below the baropause is the *barosphere*.

The terms *ionosphere* and *magnetosphere* apply, respectively, to the ionized regions of the atmosphere and to the outermost region where the geomagnetic field controls the particle motions. The outer termination of the geomagnetic field (at about ten Earth radii in the sunward direction) is the *magnetopause*.

### 1.2.2 Hydrostatic equilibrium in the atmosphere

Between them the properties temperature, pressure, density, and composition determine much of the atmosphere's behaviour. They are not independent, being related by the universal gas law which may be written in various forms, but for our purposes the form

 $P = nkT, \tag{1.1}$ 

where n is the number of molecules per unit volume, is the most useful. The quantity n is properly called the *concentration* or the *number density*, but "*density*" alone is often used when the sense is clear.

Apart from its composition, the most significant feature of the atmosphere is that the pressure and density decrease with increasing altitude. This height variation is described by the *hydrostatic equation*, sometimes called the *barometric equation*, which is easily derived from first principles. The variation of pressure with height is

$$P = P_0 \exp(-h/H), \tag{1.2}$$

5

Basic principles of the ionosphere

where P is the pressure at height h,  $P_0$  is the pressure where h=0, and H is the scale height given by

$$H = kT/(mg), \tag{1.3}$$

in which k is Boltzmann's constant, T is the absolute temperature, m is the mass of a single molecule of the atmospheric gas, and g is the acceleration due to gravity.

If T and m are constant (and any variation of g with height is neglected), H is the vertical distance over which n falls by a factor e (=2.718), and thus it serves to define the thickness of an atmosphere. H is greater, and the atmosphere thicker, if the gas is hotter or lighter. In the Earth's atmosphere H varies from about 5 km at height 80 km to 70–80 km at 500 km.

Using equation (1.1), the hydrostatic equation may be written in differential form as

$$dP/P = dn/n + dT/T = -dh/H.$$
(1.4)

From this, H can be ascribed a local value, even if it varies with height.

Another useful form is

$$P/P_0 = \exp[-(h - h_0)/H] = e^{-z},$$
(1.5)

where  $P = P_0$  at the height  $h = h_0$ , and z is the reduced height defined by

$$z = (h - h_0)/H.$$
 (1.6)

The hydrostatic equation can also be written in terms of the density ( $\rho$ ) and the number density (n). If T, g, and m are constant over at least one scale height, the equation is essentially the same in terms of P,  $\rho$ , and n, since  $n/n_0 = \rho/\rho_0 = P/P_0$ . The ratio k/m can also be replaced by R/M, where R is the gas constant and M is the relative molecular mass.

Whatever the height distribution of the atmospheric gas, its pressure  $P_0$  at height  $h_0$  is just the weight of gas above  $h_0$  in a column of unit cross-section. Hence

$$P_0 = N_{\rm T} mg = n_0 k T_0, \tag{1.7}$$

where  $N_{\rm T}$  is the total number of molecules in the column above  $h_0$ , and  $n_0$  and  $T_0$  are the concentration and the temperature at  $h_0$ . Therefore we can write

$$N_{\rm T} = n_0 k T_0 / (mg) = n_0 H_0, \tag{1.8}$$

 $H_0$  being the scale height at  $h_0$ . This equation says that, if all the atmosphere above  $h_0$  were compressed to density  $n_0$  (that already applying at  $h_0$ ), then it would

#### 1.2 Vertical structure

occupy a column extending just one scale height. Note also that the total mass of the atmosphere above unit area of the Earth's surface is equal to the surface pressure divided by g.

Although we often assume that g, the acceleration due to gravity, is a constant, in fact it varies with altitude as  $g(h) \propto 1/(R_{\rm E}+h)^2$ , where  $R_{\rm E}$  is the radius of the Earth. The effect of changing gravity may be taken into account by defining a *geopotential height* 

$$h^* = R_{\rm E} h / (R_{\rm E} + h). \tag{1.9}$$

A molecule at height *h* over the spherical Earth has the same potential energy as one at height  $h^*$  over a hypothetical flat Earth having gravitational acceleration g(0).

Within the homosphere, where the atmosphere is well mixed, the mean relative molecular mass determines the scale height and the variation of pressure with height. In the heterosphere, the partial pressure of each constituent is determined by the relative molecular mass of that species. Each species takes up its own distribution, and the total pressure of the atmosphere is the sum of the partial pressures in accordance with Dalton's law.

### 1.2.3 The exosphere

In discussing the atmosphere in terms of the hydrostatic equation we are treating the atmosphere as a compressible fluid whose temperature, pressure, and density are related by the gas law. This is valid only if there are sufficient collisions between the gas molecules for a Maxwellian velocity distribution to be established. As the pressure decreases with increasing height so does the collision frequency, and at about 600 km the distance traveled by a typical molecule between collisions, the *mean free path*, becomes equal to the scale height. At this level and above we have to regard the atmosphere in a different way, not as a fluid but as an assembly of individual molecules or atoms, each following its own trajectory in the Earth's gravitational field. This region is called the *exosphere*.

While the hydrostatic equation is strictly valid only in the barosphere, it has been shown that the same form may still be used if the velocity distribution is Maxwellian. This is true to some degree in the exosphere, and the use of the hydrostatic equation is commonly extended to 1500–2000 km, at least as an approximation. However, this liberty may not be taken if there is significant loss of gas from the atmosphere, since more of the faster molecules will be lost and the velocity distribution of those remaining will be altered thereby. The lighter gases, helium and hydrogen, are affected most.

The rate at which gas molecules escape from the gravitational field in the exosphere depends on their vertical speed. Equating the kinetic and potential energies of an upward-moving particle, its escape velocity  $(v_e)$  is given by

Basic principles of the ionosphere

$$v_{\rm e}^2 = 2gr,$$
 (1.10)

where *r* is the distance of the particle from the center of the Earth. (At the Earth's surface the escape velocity is  $11.2 \text{ km s}^{-1}$ , irrespective of the mass of the particle.)

By kinetic theory the root mean square (r.m.s.) thermal speed of gas molecules  $(\overline{v^2})$  depends on their mass and temperature, and, for speeds in one direction, i.e. vertical,

$$m\overline{v^2}/2 = 3kT/2.$$
 (1.11)

Thus, corresponding to an escape velocity  $(v_e)$  there can be defined an *escape temperature*  $(T_e)$ .

 $T_{\rm e}$  is 84000 K for atomic oxygen, 21000 K for helium, but only 5200 K for atomic hydrogen. At 1000–2000 K, exospheric temperatures are smaller than these escape temperatures, and loss of gas, if any, will be mainly at the high-speed end of the velocity distribution. In fact, the loss is insignificant for O, slight for He, but significant for H. Detailed computations show that the resulting vertical distribution of H departs significantly from the hydrostatic at distances more than one Earth radius above the surface, but for He the departure is small.

## 1.2.4 The temperature profile of the neutral atmosphere

The atmosphere's temperature profile results from the balance amongst sources of heat, loss processes, and transport mechanisms. The total picture is complicated, but the main points are as follows.

### Sources

The troposphere is heated by convection from the hot ground, but in the upper atmosphere there are four sources of heat:

- (a). Absorption of solar ultra-violet and X-ray radiation, causing photodissociation, ionization, and consequent reactions that liberate heat;
- (b). Energetic charged particles entering the upper atmosphere from the magnetosphere;
- (c). Joule heating by ionospheric electric currents; and
- (d). Dissipation of tidal motions and gravity waves by turbulence and molecular viscosity.

Generally speaking, the first source (a) is the most important, though (b) and (c) are also important at high latitude. Most solar radiation of wavelength less than 180 nm is absorbed by  $N_2$ ,  $O_2$  and O. Photons that dissociate or ionize molecules or atoms generally have more energy than that needed for the reaction, and the excess appears as kinetic energy of the reaction products. A newly created photoelectron, for example, may have between 1 and 100 eV of kinetic energy, which

#### 1.2 Vertical structure

subsequently becomes distributed throughout the medium by interactions between the particles (optical, electronic, vibrational, or rotational excitation, or elastic collisions, depending on the energy.) Elastic collisions redistribute energy less than 2 eV, and, since this process operates mainly between electrons, these remain hotter than the ions. Some energy is reradiated, but on average about half goes into local heating. It can generally be assumed that in the ionosphere the rate of heating in a given region is proportional to the ionization rate.

The temperature profile (Figure 1.2) can be explained as follows. The maximum at the stratopause is due to the absorption of 200–300 nm (2000–3000 Å) radiation by ozone ( $O_3$ ) over the height range 20–50 km. Some 18 W m<sup>-2</sup> is absorbed in the ozone layer. Molecular oxygen ( $O_2$ ), which is relatively abundant up to 95 km, absorbs radiation between 102.7 and 175 nm, much of this energy being used to dissociate  $O_2$  to atomic oxygen (O). This contribution amounts to some 30 mW m<sup>-2</sup>. Radiation of wavelengths shorter than 102.7 nm, which is the ionization limit for  $O_2$  (See Table 1.1 of Section 1.4.1), is absorbed to ionize the major atmospheric gases  $O_2$ , O, and  $N_2$  over the approximate height range 95–250 km, and this is what heats the thermosphere. Though the amount absorbed is only about 3 mW m<sup>-2</sup> at solar minimum (more at solar maximum), a small amount of heat may raise the temperature considerably at great height because the air density is small. Indeed, at the greater altitudes the heating rate and the specific heat are both proportional to the gas concentration, and then the rate of increase in temperature is actually independent of height.

At high latitude, heating associated with the aurora – items (b) and (c) – is important during storms. Joule heating by electric currents is greatest at 115-130 km. Auroral electrons heat the atmosphere mainly between 100 and 130 km.

#### Losses

The principal mechanism of heat loss from the upper atmosphere is radiation, particularly in the infra-red. Emission by oxygen at 63  $\mu$ m is important, as are spectral bands of the radical OH and the visible airglow from oxygen and nitrogen. The mesosphere is cooled by radiation from CO<sub>2</sub> at 15  $\mu$ m and from ozone at 9.6  $\mu$ m, though during the long days of the polar summer the net effect can be heating instead of cooling.

### Transport

The thermal balance and temperature profile of the upper atmosphere are also affected by processes of heat transport. At various levels conduction, convection, and radiation all come into play.

Radiation is the most efficient process at the lowest levels, and the atmosphere is in radiative equilibrium between 30 and 90 km. *Eddy diffusion*, or convection, also operates below the turbopause (at about 100 km), and allows heat to be carried down into the mesosphere from the thermosphere. This flow represents a major loss of heat from the thermosphere but is a minor source for the mesosphere.

10 Basic principles of the ionosphere

In the thermosphere (above 150 km) thermal conduction is efficient because of the low pressure and the presence of free electrons. The large thermal conductivity ensures that the thermosphere is isothermal above 300 or 400 km, though the thermospheric temperature varies greatly from time to time. *Chemical transport* of heat occurs when an ionized or dissociated species is created in one place and recombines in another. The mesosphere is heated in part by the recombination of atomic oxygen created at a higher level. There can also be horizontal heat transport by large-scale winds, which can affect the horizonal distribution of temperature in the thermosphere.

The balance amongst these various processes produces an atmosphere with two hot regions, one at the stratopause and one in the thermosphere. The thermospheric temperature, in particular, undergoes strong variations daily and with the sunspot cycle, both due to the changing intensity of solar radiation.

### 1.2.5 Composition

The upper atmosphere is composed of various major and minor species. The former are the familiar oxygen and nitrogen in molecular or atomic forms, or helium and hydrogen at the greater heights. The minor constituents are other molecules that may be present as no more than mere traces, but in some cases they can exert an influence far beyond their numbers.

### Major species

The constant mixing within the turbosphere results in an almost constant proportion of major species up to 100 km, essentially the mixture as at ground-level called "air", although complete uniformity cannot be maintained if there are sources and sinks for particular species. Molecular oxygen is dissociated to atomic oxygen by ultra-violet radiation between 102.7 and 175.9 nm:

$$O_2 + h\nu \to O + O, \tag{1.12}$$

where  $h\nu$  is a quantum of radiation. An increasing amount of O appears above 90 km. The atomic and molecular forms are present in equal concentrations at about 125 km, and above that the atomic form increasingly dominates. Nitrogen is not directly dissociated to the atomic form in the atmosphere, though it does appear as a product of other reactions.

Above the turbopause mixing is less important than diffusion, and then each component takes an individual scale height depending on its relative atomic or molecular mass (H=kT/(mg)). Because the scale heights of the common gases vary over a wide range – H = 1, He = 4, O = 16, N<sub>2</sub> = 28, O<sub>2</sub> = 32 – the relative composition of the thermosphere is a marked function of height, the lighter gases becoming progressively more abundant as illustrated in Figure 1.3. Atomic oxygen dominates at a height of several hundred kilometers. Above that is the