

**NOTICE: WARNING
CONCERNING COPYRIGHT RESTRICTIONS**



- The copyright law of the United States (Title 17, United States Code) governs the making of photocopies or other reproductions of copyrighted material.
- Under certain conditions specified in the law, libraries and archives are authorized to furnish a photocopy or other reproduction. One of these specific "fair use" conditions is that the photocopy or reproduction is not to be "used for any purpose other than private study, scholarship, or research." If a user makes a request for, or later uses, a photocopy or reproduction for purposes in excess of "fair use," that user may be liable for copyright infringement.

Exploring Geographic Information Systems

Second Edition

Nicholas Chrisman

University of Washington



John Wiley & Sons, Inc.

New York Chichester Brisbane Toronto Singapore Weinheim

ACQUISITIONS EDITOR Ryan Flahive
DEVELOPMENT EDITOR Joan Petrokofsky
SENIOR PRODUCTION EDITOR Elizabeth Swain
MARKETING MANAGER Clay Stone
SENIOR DESIGNER Dawn Stanley
COVER DESIGNER Harold Nolan
ILLUSTRATION EDITOR Sandra Rigby
COVER ILLUSTRATION "Automated Radio Design Support-Teligent RF Engineering" by Jubal Harpste, Mike Ruth, and Brian Sandrik. Selected graphic image supplied courtesy of Teligent IT/Applications and Environmental Systems Research Institute, Inc. Copyright © 1999, Teligent IT/Applications.

This book was set in 10/12 New Caledonia by York Graphic Services, Inc.-Shippensburg Facility (TechBooks) and printed and bound by Malloy. The cover was printed by Lehigh.

This book is printed on acid-free paper. ☺

Copyright ©2002 by John Wiley & Sons, Inc. All rights reserved.

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 605 Third Avenue, New York, NY 10158-0012, (212) 850-6011, fax (212) 850-6008, E-Mail: PERMREQ@WILEY.COM.

To order books or for customer service, call 1(800)-CALL-WILEY (225-5945).

Library of Congress Cataloging in Publication Data:

Chrisman, Nicholas R.

Exploring geographic information systems / by Nicholas Chrisman.—2nd ed.

p. cm.

Includes bibliographical references (p.).

ISBN 0-471-31425-0 (pbk. : alk. paper)

1. Geographic information systems. I. Title

G70.212.C48 2002

025.16'91—dc21

2001017548

Printed in the United States of America

10 9 8 7 6 5 4 3 2

Defining a Geographic Information System

There are dozens of definitions for the term *geographic information system* (GIS), each developed from a different perspective or disciplinary origin. Some focus on the map connection; some stress the database or the software tool kit; and others emphasize applications such as decision support (Maguire 1991; Chrisman 1999a). One of the most general definitions was developed by consensus among 30 specialists:

Geographic Information System—A system of hardware, software, data, people, organizations and institutional arrangements for collecting, storing, analyzing and disseminating information about areas of the earth. (Dueker and Kjerne 1989, pp. 7–8)

While this definition may seem bland, it encompasses all the characteristics, as long as the terms are expanded to their intended meaning. For example, the word *system* implies a group of connected entities and activities. An automated information system organizes a collection of data, computer procedures, and human organizations to serve some particular purpose. For a GIS, the purpose could involve a complex decision, such as the policy for timber harvest, or a routine decision, such as granting a permit or maintaining an inventory. Notice that the definition carefully distinguishes between the data in the system and the information that results from the system. Data provide the raw material for information, much as map symbols convey a map message. For both maps and information systems, the raw data are not enough; additional relationships must be constructed from the context.

The most common understanding of a GIS emphasizes that a GIS is a tool. However, no tool is totally neutral; a GIS can be designed to be effective and efficient for

a certain purpose. Tools are developed within a social and historical context to serve changing needs, but tools are also intended to change their environment—as the story of the project in Dane County demonstrates. The perspective of this book can be summarized by the following definition:

Geographic Information System (GIS)

The organized activity by which people

- measure aspects of geographic phenomena and processes;
- represent these measurements, usually in the form of a computer database, to emphasize spatial themes, entities, and relationships;
- operate upon these representations to produce more measurements and to discover new relationships by integrating disparate sources; and
- transform these representations to conform to other frameworks of entities and relationships.

These activities reflect the larger context (institutions and cultures) in which these people carry out their work. In turn, the GIS may influence these structures.

The rest of this book will follow the sequence of this definition, exploring how a GIS works to provide solutions for problems.

REFERENCE SYSTEMS FOR MEASUREMENT

CHAPTER OVERVIEW

- Introduce reference systems for time, space, and attributes.
- Extend Stevens' levels of measurement to provide a richer basis for attribute measurement.

HOW INFORMATION WORKS

An exploration of geographic information systems must begin with an operational understanding of information, then tighten the focus on what is specifically geographic, and finally examine the system components. The word *information* appears nearly everywhere in current life. We are said to be in an Age of Information, much as earlier and not so distant periods were termed the Age of the Airplane and the Age of the Atom. Such trendy catchphrases reflect what is economically important at a given time, but the glare of the spotlight may flatten out important details. **Information** occupies a middle stage in a process modeled on the scientific method. The starting point involves data—raw observations, that have no particular value by themselves. Somehow, through procedures not often totally explained, these raw data acquire value when placed in a frame of reference—a system of relationships among objects and assumptions about those relationships. For example, the digits 8.5 do not mean much unless you know they measure a water level in meters vertically from extreme low tide. This process also implies that information leads to higher levels of knowledge, through further refinement and interpretation.

Information: Data (observations, measurements, etc.) placed in context of a system of meaning (a set of relationships and assumptions about those relationships). Information, built into larger context, constructs knowledge.

This sequential flow from data to information and eventually to knowledge does not occur without human beings actively engaged. Gigabytes of databases do not produce refined information products unless someone does a lot of work. In addition, there is no guarantee that everyone will extract exactly the same results when confronted with the same data. Knowledge is not simply a passive result of assembling the data. It takes special talent to be willing to abandon your assumptions and to be open to surprise.

This exploration will not reveal a specific boundary where the data instantly become useful information. The whole book is about the process of providing context and of discovering relationships. Constructing information requires experimentation and exploration. It is not some smooth, guaranteed progression as on an assembly line. Understanding often comes with a spark of recognition—a realization that some particular fact jolted you to recognize a new relationship. Though this moment may seem distinctly personal, the practice of measurement provides guidance from the shared experience of communities of scientists, government officials, and business people. Measurement techniques ensure that separate facts relate to a common reference system. Geographic information depends on common forms of measurement, although it raises some particular issues and special problems.

BASIC COMPONENTS OF GEOGRAPHIC INFORMATION

Geographic information is commonly broken into the components of *space*, *time*, and *attribute*. Space, although it is an obvious component of geographic information, can be understood from a number of different perspectives. At the simplest, a space can consist of distinct “places” that are only different from each other. However, it does not require much observation to figure out that certain places are nearer than others and that some align in the same direction. Building up these relationships leads to the basics of geometry. The world of sensory experience is basically three dimensional. Objects have length, width, and height, and each is located at some distance and direction from the others. When dealing with larger geographic regions, attention can be limited to a thin shell of the earth’s surface. For much mapping (and GIS), this space is predominantly two dimensional. Even though it is often convenient to use **Cartesian** geometry within confined regions, the surface of the earth is not a limitless plane. The most common spatial reference systems for GIS provide a translation between the three-dimensional shape of the earth and the flat spaces traditionally used for maps.

Time often plays a silent role in maps, though there is always some implicit or explicit temporal reference. The most common map works like a snapshot—valid for a specific moment in time. Of course, time is more than a collection of points. The inexorable progression of time and the inability to turn back the clock lead to the most

Cartesian: Geometry of unbounded spaces, particularly in two dimensions; the geometry of a flat plane treated by the tools of analytical geometry (cartesian coordinates) developed by René Descartes in the seventeenth century.

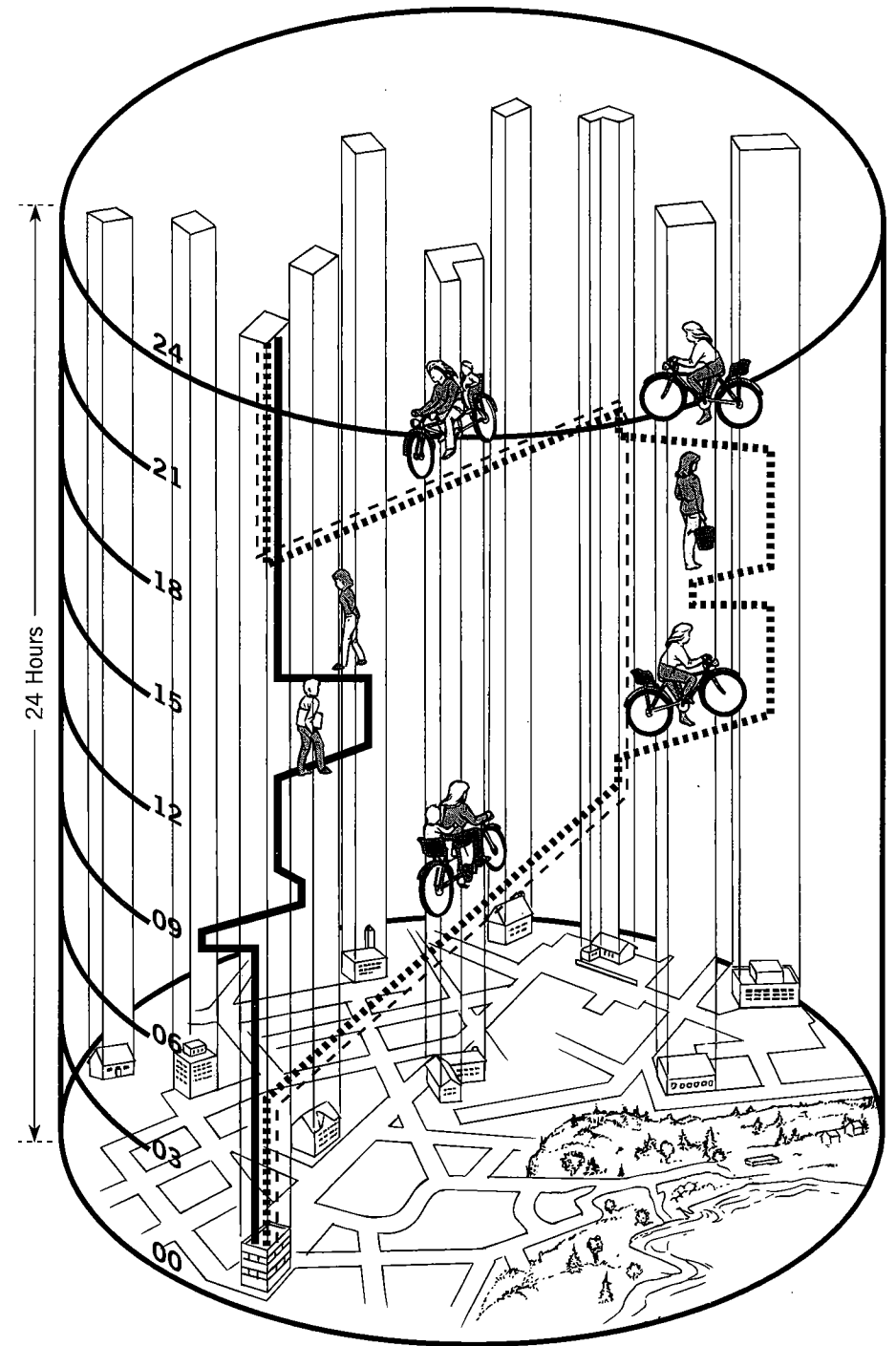


Figure 1-1 Space-time diagram showing people in their daily lives. Redrawn from Parkes and Thrift (1980).

prevalent view of time as an arrow—a line with one direction. Sometimes, however, it makes sense to treat time as cyclical, such as when astronomical and climatic events repeat. For example, when the earth returns to a particular position in its orbit, the stars return to the positions observed the previous year. In these cases, the sense of cycles comes from the connection between the spatial realm and the temporal, not because any span of time actually repeats.

The interplay of time and space in daily life can be read from the time-space diagrams produced by time geographers [following the lead of Hägerstrand (1970)]. Figure 1-1 shows a day (0 to 24 hours vertically) in a small town. The buildings remain in place, appearing as pillars (unless a new one appears or an old one gets demolished on this particular day). Members of one family start their day together, their movements charted as diagonal lines starting from their apartment building. Quite early, the mother takes off with the younger child on her bicycle, leaves the child at a day-care center (light dashed line), and continues on with her daily activities (thick dashed line) until she returns to the day-care center and then home. The other child (dark line) stays closer to home, going to school and to some afternoon lessons, returning home before the mother returns. The social interaction of this town is strongly synchronized by clocks. The mother's life is a dance of social expectations to be in certain places at particular times. The movement between places and times creates the social web of the town. This diagram illustrates how time and space provide a frame of reference for all activities.

The third component of geographic information, the **attribute**, can range from observable physical properties to aesthetic judgments. In Figure 1-1, there are dozens of possible attributes, including the width of streets, the names of the people, and the scenic beauty from each viewpoint. Information extracted from the time and space dimensions, such as a rate of speed, can also be treated as attributes. The attributes attached to a GIS are certainly the most varied of the three components. While this chapter can give technical details for the measurement of time and space, it can suggest only general classes for attributes.

Reference Systems

Each of the three components of geographic information is measured with respect to some particular **reference system**. Such a system provides rules to interpret individual observations with respect to others and to document the rules so that results can be repeated and compared. The technology of temporal reference systems is quite ancient, since calendars and clocks have existed for millennia. The other forms of reference systems are more recently established and less universal. The clearest reference systems to use are those established by explicit standards, though workable results can come from less formalized procedures as long as they are shared by all users.

Attribute: The range of possible values of a characteristic; an attribute value is a specific instance of the characteristic associated with a geographic feature.

Reference system: An established set of rules for measurement. Provides a means to compare a particular measurement to others performed with reference to the same set of rules. Geographic information requires reference systems for time, space, and attributes.

Temporal Reference Systems Time, with its strong sense of a linear order, is simpler than space. The linear axis of time—measured in units such as seconds, hours, and years—orders our lives in many ways (Figure 1-2). A simple temporal reference system merely requires an origin (a time to call zero) and a unit of measurement, such as a second. A stopwatch starts from its own zero each time you start it. A more complex system, such as a calendar, requires rules for counting days in months, and as long as others follow the same rules it remains reliable. Living together in a civilized society creates the need for a shared temporal reference system. Each ancient civilization created its own calendar, some based on the lunar cycle, others based on the solar year. During the past century, global activity became synchronized by a common reference time (Greenwich Mean Time) and a common calendar. Each time zone around the world sets its clocks so that solar noon corresponds roughly to 12 o'clock. A time zone is a subsidiary reference system based on an offset from Greenwich Mean Time. Similarly, certain countries and religious communities retain alternative calendars, but the correspondence is well understood. By adopting a common reference system, time measurements can be compared and mathematical operations such as subtraction produce useful results.

Some aspects of time are cyclical. In environmental studies of all kinds, the seasons play an important role. A “growing year” can be thought to start in spring while a “water year” might start in fall. There is no necessary starting point to a yearly cycle.

Another kind of reference system consists of ordered periods. For example, administrative procedures often specify a sequence of events, perhaps with some guidelines for duration, without respect to any particular starting point. Thus, an environmental impact statement might have a series of planned phases such as a scoping process, public comment, the analysis phase, a draft report, more public comment, then a final report. Any particular project can be located along this time sequence without needing a numerical measure.

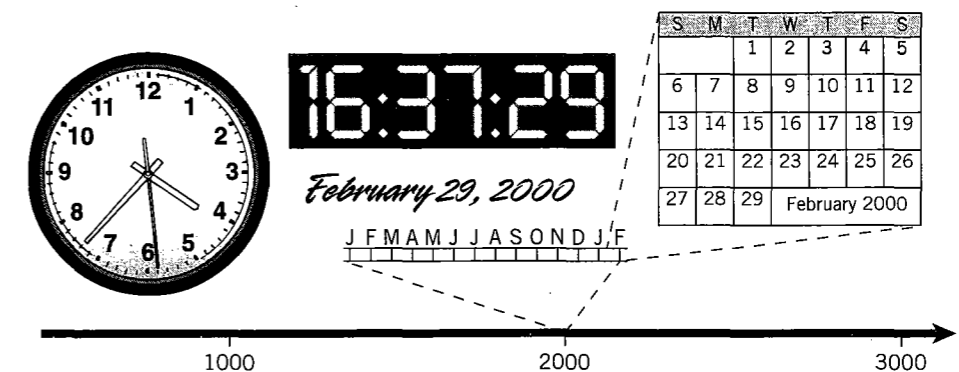


Figure 1-2 Temporal reference systems, lines and cycles. Clocks with rotary hands or digital displays repeat a cycle of measurements, while a time line continues indefinitely. The conventional calendar is a complex combination of cyclical elements at different frequencies (weeks, months, years, leap years, etc.).

Spatial Reference Systems Perhaps the family members in Figure 1-1 navigate through their town as a network of named places: buildings, roads, parks, and neighborhoods. They may never see it from above or draw a map of it, but still they would eventually learn certain geometric rules (such as the triangular inequality—the straight line home is shorter than combining it with an errand). The relationships between places create a geometry, even if it is not formalized and measured. A **spatial reference system** provides a refined tool by formalizing these relationships using analytical geometry.

Beyond specifying an origin and a unit of measure, spatial measurement requires geometric assumptions because an object cannot be located in a two-dimensional space by a single number along one axis. Spatial measurement for a limited site can use a local grid laid out in the field as long as the two axes are strictly laid out at right angles. Here the geometric model is a flat plane. The spatial reference for the La Selva Biosphere Reserve project introduced earlier was about as simple as it could be. Like many scientific sites, the researchers are more interested in local coordinates and the ability to find things in the field, not the relationships to places outside the Reserve. At La Selva during the 1980s, ecologists placed pieces of plastic pipe in the ground on a 200-meter grid. The idea was that researchers could locate themselves using a compass and a tape measure relative to the nearest pipe. The map the students had to digitize was created with this 200-meter grid as its reference system.

In 1991, the Reserve adopted a new spatial reference system and installed another set of pipes spaced 50 meters apart along one axis and 100 meters on the other. This new reference system provided a much denser grid of plastic pipes in the field and was laid out using improved surveying instruments. Due to imperfections in the original survey, the 200-meter grid was not on exactly the same axes as the successor; it was rotated at a small angle. The 1980s map could not be incorporated directly into the Reserve's GIS because of the difference in the spatial reference system. Because both systems were totally local, the reference systems could only be related by measuring some objects in both systems to derive the geometric relationships between the two planes. A student did this work in the field by measuring the markers of the 1980s system according to the 1991 reference system.

Local coordinate systems disconnected from a geodetic framework will probably die out as positioning technology becomes more common and the demand for GIS integration increases. La Selva with its deep "cloud forest" canopy may be one of the last places to make the conversion. In any case, GIS workers who must salvage records from the past will have to cope with obsolete systems for a long time to come.

While some local mapping projects can still be performed using an isolated, planar reference system, a GIS must usually mobilize a more complex process with a series of geometric steps that mobilize the science of **geodesy** [see Defense Mapping

Spatial reference system: A mechanism to situate measurements on a geometric body, such as the earth; establishes a point of origin, orientation of reference axes, and geometric meaning of measurements, as well as units of measure.

Geodesy: Science of measuring the shape of the earth and establishing positions on it. Involves study of geophysical properties such as variations in gravitational field. Adjective form: geodetic.

Agency (1984) for a comprehensive review] and analytical cartography. The actual shape of the earth, the **geoid**, is too lumpy to use as a reference surface. The first step adopts a model of the earth, usually in the form of a reference **ellipsoid**. There are dozens of ellipsoids in use, each chosen to fit the apparent shape of the earth in the regions surveyed. Each of these reflects the transitional period when local surveys were connected together in larger networks, but, until recently, it was technically challenging to connect all these networks into a global system. Once connected through global geodesy, the variations since the 1970s have become effectively insignificant for mapping purposes (Table 1-1).

While the ellipsoid provides a smooth surface, it is featureless and thus not sufficient as a reference system. A **geodetic datum** populates an ellipsoid with specific points whose locations have been established through astronomical surveying and careful adjustment to compensate for errors. It is a long-standing practice to specify positions on an ellipsoid using coordinates of latitude and longitude given as angles (Figure 1-3a). This geometric model uses the axis of the earth's rotation as its north-south axis. The plane of the equator—at right angles to this axis—provides the origin for angles north and south (latitude). Longitude measures the angles east and

TABLE 1-1 World Geodetic Standards: Reference Ellipsoids^a

<u>Name</u>	<u>Equatorial (Major) Axis in Meters</u>	<u>Flattening (1/f)</u>	<u>Region</u>
Airy 1830	6377563	299.325	Great Britain
Bessel 1841	6377397.2	299.153	Central Europe
Everest 1830	6377276.3	300.80	Indian subcontinent
Clarke's 1866	6378206.4	294.98	North America
Clarke's 1880	6378249.2	293.47	Africa; France
Krasovsky 1940	6378245	298.2	Former Soviet Union
World Geodetic System 1972	6378135	298.26	NASA, U.S. military
GRS 1980/ WGS 84 ^b	6378137	298.257	GPS, new systems

^aThese reference ellipsoids may serve as the best fit to the actual geoid in different parts of the world. A horizontal "datum" adopts a reference ellipsoid and locates geodetically surveyed points on that ellipsoid.

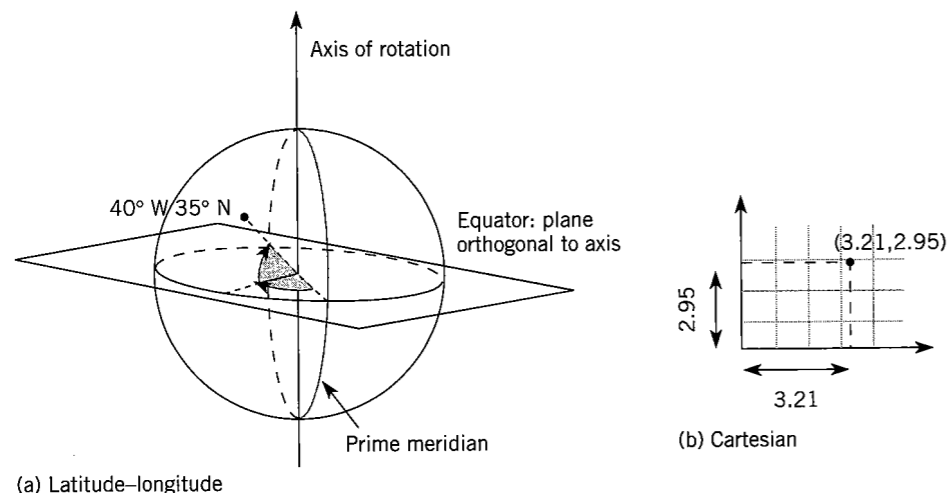
^bAt the resolution shown, Geodetic Reference System 80 and World Geodetic System 84 are the same.

Source: Snyder (1987), Table 1, p. 12.

Geoid: Three-dimensional shape of the earth defined by the surface where gravity has the value associated with mean sea level.

Ellipsoid: Three-dimensional object formed by rotating an ellipse around its minor axis; an oblate ellipsoid approximates the shape of the earth (geoid), computed by the best fit to geodetic observations. See Table 1-1.

Geodetic datum: A geodetic reference system, usually divided into vertical and horizontal standards. A horizontal datum is based on a given ellipsoid and specified latitude-longitude coordinates for certain points. The plural of datum is *datums*, despite the word's Latin origins. (Not to be confused with the usage of *datum* as a single value, plural *data*.)



(a) Latitude-longitude

Figure 1-3 Spatial reference systems: Latitude-longitude pairs measure angles from the plane of the equator (orthogonal to the rotational axis) and the prime meridian (the meridian of Greenwich, England). The World Geodetic Reference System specifies a particular ellipsoid that provides the size of the earth to use with this model of angular measurement. Other reference systems may use a plane or a cone as a projection surface, connected to the geodetic reference system. Measurements on these flat projection surfaces are usually distances from a pair of axes (cartesian coordinates).

west around the equator using an arbitrary origin set by international convention along the meridian through Greenwich, England.

A particular place on the earth will have different values of latitude-longitude, depending on the choice of datum and ellipsoid. One common source of incompatibility between spatial reference systems in the United States comes from the transition between the **North American Datums (NAD)** of 1927 and 1983. This includes a change in reference ellipsoids plus the more local effects of readjusting the geodetic survey data. For example, the first-order **geodetic control** station Northwest Blake (PID SY5304) has been located on the same boulder on the beach of Blake Island in Puget Sound since it was first established in 1857. The coordinates assigned to this point have changed due to surveying technology and the development of more comprehensive geodetic models. Using the most current geodetic adjustments, this marker has a coordinate 4.46 seconds of longitude farther

North American Datum (NAD): An adjustment of geodetic measurements that provides the accepted horizontal reference system for North America. The 1927 Datum held Mead's Ranch, Kansas, as a fixed point, while the 1983 Datum performed a simultaneous adjustment of all measurements. NAD 1927 uses Clarke's 1866 ellipsoid, while NAD 1983 uses the 1980 Geodetic Reference System.

Geodetic control: Reference marks used to establish a spatial reference system for a specific project; under ideal circumstances these are geodetic survey points established as part of the National Geodetic Reference Network.

west of Greenwich than it had using NAD27. Such a difference means about one kilometer at this latitude, although the marker has not moved at all. For this reason, knowing the full spatial reference system can be very critical when registering sources compiled at different dates. Since ellipsoids are mathematically derived, they can be converted from one to another, but two datums require further measurements for accurate conversion. Similar differences occur between some adjacent European countries.

While geodetic reference systems are most appropriate for geographic information, maps are not constructed on ellipsoidal geometry. Another step is required. Most maps adopt a simpler geometric model (such as a cylinder or a cone) positioned with respect to the geodetic model—a process called **projection**. This technique was developed to permit map construction on two-dimensional media, but the simplicity of cartesian coordinates has been retained for many software packages. Cartesian coordinates measure distances along two axes at right angles (Figure 1-3b), but when defined through a projection, these positions retain their link to the global geodetic reference system. There are hundreds of possible projections and infinite ways to position them on the ellipsoid (Snyder 1987). In practice, for a given area a small number of projections will be used. Projections can be classified by the properties they preserve as well as by their geometry (whether based on a **developable surface** such as a cone or cylinder).

Selection of a map projection is often a matter of tradition, linked to the historical development of a given set of map users. Nautical charts, for example, are nearly always drawn on the Mercator projection, with the equator as line of tangency. This cylindrical projection is **conformal**, and it has a special property, namely, that corrected compass bearings are straight lines. Manual plotting of courses and positions requires these particular properties. For topographic mapping, most of the world is covered at some scale in either a **transverse** Mercator such as Universal Transverse Mercator (**UTM**) or a conformal conic such as Lambert Conformal Conic (Table 1-2). This choice was based on the needs of artillery officers and engineers, who originally developed these maps. Equal-area projections might be more useful for many kinds of GIS analysis.

Projection: Geometric transformation that converts latitude-longitude coordinates into planar coordinates. Projections can be based on a developable surface (such as a plane, cylinder, or cone) or on a mathematical function.

Developable surface: Three-dimensional object that can be flattened into a plane without scale distortions. Cylinders and cones, since they curve in only one axis, can be converted into planes by making a single cut. Common projections use developable surfaces but some projections use more complex functions.

Conformal: Property of a projection that preserves the shape of geographic features. Within an immediate vicinity of a point, angles are preserved, but linear and areal scales have to be sacrificed to obtain this property.

Transverse: A projection oriented at right angles to the equator. A transverse cylindrical projection uses a meridian of longitude as its central meridian.

UTM: Universal Transverse Mercator; a spatial reference system using a set of transverse Mercator projections 6° wide that cover the earth (except for polar regions covered by two polar stereographic projections).

TABLE 1-2 Common Projections Used as Spatial Reference Systems

Transverse Mercator Systems

Worldwide	Universal Transverse Mercator (6° strips)
United Kingdom	National Grid
Germany	Gauss-Kruger (3° strips)
State Plane System ^a (USA)	North-south states, e.g., Illinois in two zones

Lambert Conformal Conics

France	Grille Lambert, three zones
State Plane System ^a (USA)	East-west states, e.g., Washington in two zones

Other Projections

Malaysia (peninsular)	Malayan skew orthomorphic (oblique Mercator)
Switzerland	National Grid (Gauss-Kruger)
State Plane System ^a (USA)	Alaskan Panhandle; oblique Mercator

^aThe State Plane System for the United States defines 125 zones, some as small as single counties.

A common system of spatial reference is a critical element of a GIS, since it brings different map layers into correspondence (Figure I-1). The Dane County project (see Introduction) discovered a half-dozen projections in active use, plus many maps without any known projection (much like the La Selva local coordinate systems). The project performed geodetic surveying to relate the different reference systems, then transformed each source following procedures described in Chapter 3.

Some countries try to limit incompatibilities by ordaining a particular projection as a common spatial reference system. All maps in Switzerland and the United Kingdom are referenced to their respective national grid systems. Larger countries and most American states cannot cover their territory with a single projection zone without distortions beyond the tolerances required for many applications. Most projects require translation between multiple spatial reference systems.

With computer representation, converting between geodetically referenced systems poses no real difficulty. A map in one projection can be converted to latitude-longitude (through the inverse of the projection function) then projected into some other datum and projection. Some calculation time may be required, but on current computers it is no great burden. The U.S. government has placed the General Cartographic Transformation Package (GCTP) in the public domain, and many commercial packages have incorporated this software.

Attribute Reference Systems Just as reference systems apply to time and space, they also apply to attributes. Unlike the common approaches that apply to time and space, each particular attribute scale requires its own reference system. There are some general rules for attribute measurement that are widely used in cartography and social science statistics. The next section reviews some of these concepts, with some additional rules related to common geographic attributes.

LEVELS OF MEASUREMENT

The wide diversity of attributes indicates a huge number of techniques for measurement. To a purist, like a classical physicist, measurement provides a numerical relationship between some standard object and the object being measured. Consider the attribute *length*. Every entity in space can be measured by comparing its length to some other length. The procedure begins by placing a standard measuring rod alongside the object to be measured, marking where the end of the rod falls, then placing the rod again beginning at the mark, until the end of the object is reached. This procedure implements a physical form of addition. The number of times the rod is placed represents the ratio of the length of the object to the length of the rod. Using similar comparisons, physicists developed procedures to measure temperature, mass, electrical charge, and more.

In nineteenth-century physics, fundamental physical properties were considered *extensive* because they *extended* in some way as length does in space. Other properties, like density, were built up as ratios of the extensive properties and were thus *derived*. The fundamental physical properties form the basis for the international standards comprising the metric system or **SI**. But the attributes used for geographic information reach far beyond the SI measures and ratios derived from them.

To provide a framework for a broader range of measurement types, Stanley Stevens (1946), a psychologist at Harvard University, published an article in *Science* proposing a framework based on what he called **levels of measurement**. Stevens adopted a very simplified definition of measurement as the “assignment of numbers to objects according to a rule” (Stevens 1946, p. 677). Stevens’ schema has become a basis for social science methods and a framework for cartography and GIS. Because Stevens’ classification is often misapplied and misinterpreted, the levels of measurement deserve careful scrutiny. Some revisions and extensions must be considered to accommodate geographic information.

Stevens used the concept of **invariance** under transformations. Invariance considers the degree that a **scale** can retain its essential information content even if it is not identical to some other scale. A level groups the scales that share a set of possible transformations. One example of invariance involves temperature. A measurement in °F can be transformed into °C without loss of information; thus, these mea-

SI: Système International d’Unités; the system of weights and measures established by international agreement in 1875. The International Bureau of Weights and Measures in Sèvres, France, oversees the measurement standards. SI defines seven base units from which many others can be derived: meter for length, kilogram for mass, second for time, kelvin for temperature, ampere for electric current, mole for chemical quantity, and candela for intensity of light.

Level of measurement: A grouping of measurement scales based on the invariance to transformations. A measurement scale at a given level of measurement can be transformed into another scale at the same level without losing information.

Invariance: Properties that remain unchanged despite transformations of the numbers used to represent the measurement.

Scale: When applied to a scale of measurement, a system used to encode the results of a measurement; typically a number line, but generalized to include a list of categories.

measurements are at the same level. So, despite having different zero values and different units of measurement, the two scales can be related to each other. By contrast, a scale consisting of {cold, warm, hot} cannot retain all the information recorded in either Fahrenheit or Celsius scales.

The following sections explain Stevens' four levels of measurement and illustrate them with a common example. Imagine a marathon in which contestants become associated with certain attributes or measurements.

Nominal

At the most basic level, Stevens described a nominal "scale" in which objects are classified into groups. Any assignment of symbols can be used, so long as the distinct nature of each group is maintained. A nominal measure is based on set theory. The use of the word *scale* for a nominal measurement may evoke the traditional number line, but there is no such ordering implied.

In the marathon example, each contestant gets a number to wear. What does this number mean? Is it a measurement? If the number is simply pulled randomly out of a box, it has to be considered an arbitrary symbol (like a word or an icon). Other nominal attributes could be determined, such as the set of contestants wearing red shirts (Figure 1-4). Another nominal grouping might allocate contestants into either the women's event or the men's event. Any numerical symbol for these two categories (0 and 1, or 1 and 2, or 359 and 213) would be totally arbitrary. In this sense, nominal data remains invariant under the most extreme alterations; any symbol can be converted to another symbol, as long as they remain distinct from each other.

Ordinal

The ordinal level introduces the concept of an ordering. An ordinal scale applies when objects can be sorted in some manner; such a scale can exist in many forms. The most exhaustive form orders all objects completely without any ties (Figure 1-5). An example is the order of runners finishing the race (first, second, third, . . .). It makes sense to use the word *scale* for such an ordering because each successive ele-

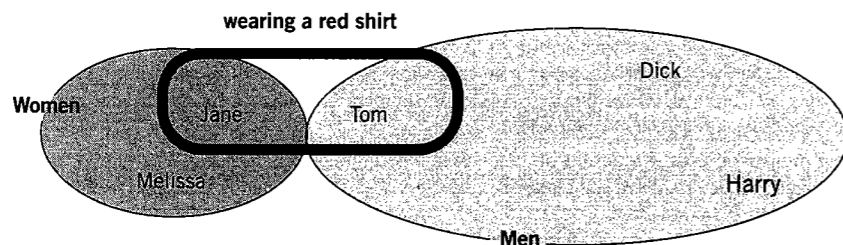


Figure 1-4 Nominal measures are not on scales at all. They create categories that can be treated as sets.

<u>Order of arrival of contestants</u>	<u>Ordinal</u>	
	<u>Women's race</u>	<u>Men's race</u>
First	Jane	Tom
Second	Melissa	Dick
Third	Leila	Harry

Figure 1-5 Strictly ordinal scales can arise from a total ordering, but ordinal scales may also arise from partial orderings.

ment follows in the same direction. In another example, the number on the shirt of each contestant in the race might also represent an ordinal measure, if the numbers are handed out of the box sequentially. Thus a seemingly nominal identifier might hide an ordering based on time or on some other property.

Not all orderings are as well behaved as this ideal model without ties. The more ties there are, the less scalelike the ordering becomes. For example, some ordinal categories use a semantic scale of words. Soils are ordered from "poorly drained" through "somewhat poorly drained" to "well drained" and "excessively drained." Opinion polls use orderings such as "strongly disagree" to "strongly agree." The ordinal level covers a wide range of possibilities. Some scales behave in a nearly numerical way, whereas others are barely evolved from a nominal level.

Ordinal measurement does not constrain numeric representation very much. It may be conventional to give out the numbers 1, 2, 3, . . . to finishers in a race, but the numbers could be any increasing sequence (0.5, 0.66, 0.75, 0.8, . . .) or (1, 3, 597, 6667, . . .) because the order is all that matters. In the example of soil drainage, we do not know if the step from "poorly drained" to "somewhat poorly drained" is identical in magnitude to the step between "well drained" and "excessively drained." For different analytical purposes, the importance of each step in the scale might vary. Some orderings may relate to an underlying numerical scale, and others may not. For example, it is hard to assume that "good" means the same thing for all respondents in an opinion poll. Ordinal values are essentially categories without the arithmetic properties usually ascribed to numbers. Hence, nominal and ordinal measures are sometimes grouped together as *categorical* measurements. It is important to remember these limitations when some GIS user wants to standardize rankings on a scale from 1 to 9. Encoding with numbers does not automatically make arithmetic valid.

Interval

In Stevens' scheme, the quantitative realm begins with interval scales that give numbers algebraic meaning. An interval scale involves a number line with an arbitrary zero point and an arbitrary interval (the unit of measurement). Thus, interval scales can be shifted by changing the zero without changing the meaning of the measurement. For example, years can be recorded on the Gregorian calendar (A.D.), the Islamic calendar (1 A.H. is A.D. 622), or the geologists' calendar [0 Before Present (B.P.)

is A.D. 1950]. In all these systems, the numerical value of a year has no particular significance. The year 2000 is not twice the year 1000 in some magnitude.

In the case of the marathon, we could assign arrival times to runners by simply noting the clock time for each arrival (Figure 1-6). As long as some basic assumptions are valid, particularly that all runners departed at the same time, then these numbers capture all the ordinal results. In addition, the differences between arrival times can be interpreted. Some of the arrivals are closer to the next arrival than others, establishing a truly numerical measure of difference between values. An elapsed running time can be twice as long as another, for example.

Ratio

Arrival times for a race provide raw results awaiting further processing. Contestants would obtain a more useful measure by subtracting the time at the start from the time at their finish. In fact, a difference between two interval measures becomes a measure on Stevens' next level: ratio.

In measurement theory, the ratio level gets the most attention. Ratio measures retain the arbitrary unit of measure from the interval scale but substitute a true origin (zero value). These properties support the arithmetic operations of addition, subtraction, multiplication, and division. On a ratio scale, if a value is twice that of another, then it represents a doubling of the quantity. The easiest ratio measures to visualize are classical extensive quantities. In the race example, the elapsed time in running the race is a ratio measure obtained by subtracting two interval measures for the start and finish (Figure 1-7). The ratio measure of elapsed running time contains all the information of the ordinal scale for ranking winners, plus it adds the numerical properties that measure how fast each contestant ran. It is clear that these ratio measures convey more information and permit more analytical treatment.

Extensive and Derived Scales

Stevens tried to combine extensive and derived measurement into one level. He defined the ratio level based on the invariance related to the arbitrary *unit of measure*. Thus, a length in feet can be converted to meters with no real change in length. But the invariance group misses some important distinctions between geographic mea-

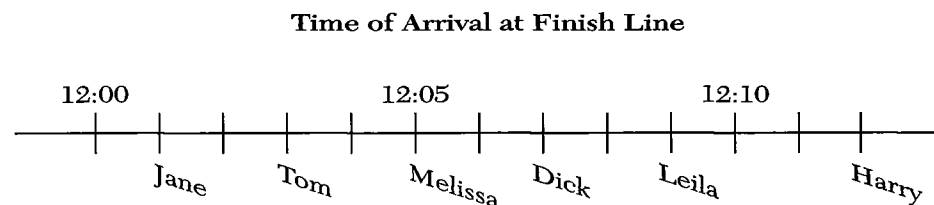


Figure 1-6 Interval scales mobilize a number line, but the origin and the unit are arbitrary.

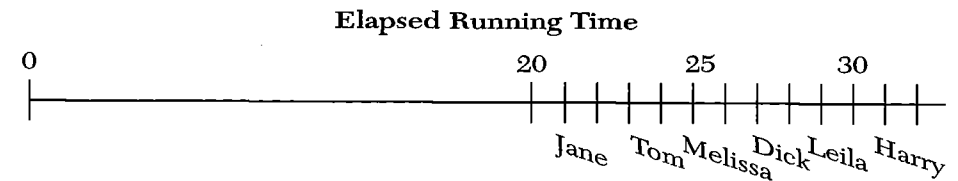


Figure 1-7 Ratio scales, the classical ideal for physical measurement, have a true origin and an arbitrary unit of measure.

surements. Consider some economic activity in a county measured in dollars. This scale is arbitrary because yen would work as well as dollars. The value of zero applies to any scale. The value for the county could be combined by addition with the value of others to obtain a measure for the state economy. Cartographic design principles would suggest **proportional symbols** for such raw values (for example, Figure 1-8 shows the total expenditure by the Department of Defense for each California county in 1976). In place of total expenditure, the county can be given a per capita measure by dividing the dollar figure by the number of persons in the county. Such a transformation removes the influence of population and concentrates on relative expenditure. This value is just as much a "ratio" as is the total expenditure figure, but cartographic rules suggest a **choropleth map** presentation for these per capita figures (Figure 1-9). Per capita values of two counties cannot really be added together because the denominators (populations) might be totally different. Notice that some counties with low total figures can have high per capita figures.

Stevens' system has been used to suggest which statistical methods apply to a given measurement. Many introductory statistics books, particularly for the social sciences, connect the levels of measurement to a group of appropriate tools. Similarly, cartography texts connect the cartographic tool kit of **graphic elements** to specific levels of measurement. The connection is not entirely straightforward, as demonstrated by the case of the California county data. For geographic data, it would make more sense to divide ratio measures into the invariance classes applied in selecting thematic map types. This would also separate those measures that are aggregated by

Proportional symbols: A thematic mapping technique that displays a quantitative attribute by varying the size of a symbol. Typically, proportional symbols use simple shapes such as circles and are scaled so that the area of the symbol is proportional to the attribute value. Each symbol is located at a point, even if it represents data collected for an area.

Choropleth map: A thematic mapping product that displays a quantitative attribute using ordinal classes applied as uniform symbolism over a whole areal feature. Sometimes extended to include any thematic map symbolized using areal objects.

Graphic elements: The characteristics of a symbol system that can be manipulated to encode information. For cartography, these include size, shape, hue, saturation, brightness, orientation, and pattern. [See Robinson et al. (1995)].

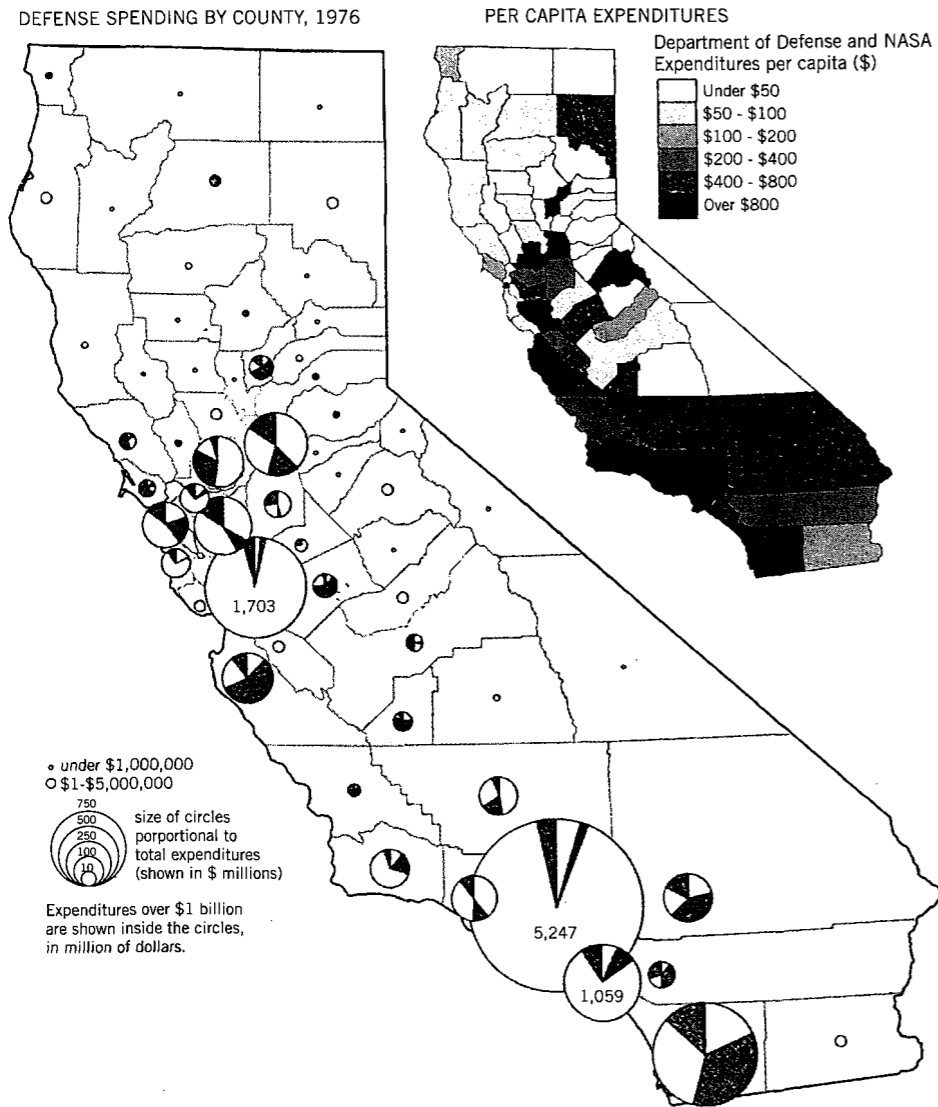


Figure 1-8 Proportional symbols use the graphic variable of size for a simple geometric symbol such as a circle. They are considered appropriate for raw measures, such as total population or total economic output. This map portrays defense spending for each county by scaling the area of the circle to be proportional to the dollar figure. (Source: Donley et al. 1979, p. 46.)

Figure 1-9 Choropleth maps use the spatial object as the symbol. The graphic variable size cannot be used without a cartogram. Here, as in many cases, the graphic variable value (a gradation from light to dark) shows the range of the attribute. Since the area of the object is a part of the symbol already, this method is most appropriate for density measures, or a derived ratio such as dollars per capita. (Source: Donley et al. 1979, p. 46.)

addition (extensive) from those that must be weighted (such as derived ratios). These are just a few examples of the operations that form the main objective of this book in later chapters.

Perhaps the best way to explain why attribute reference systems matter is by a counterexample. Along highways, it is common to announce the towns and villages through which the road passes. Most highway signs announce the name with some extra information, such as population or elevation. One town in California has a sign that takes the spirit of local pride to an extreme (Figure 1-10). Adding these three numbers is clearly a joke; yet, professionals who work with geographic information often commit equally meaningless combinations with no humorous intent. The number 4663 measures nothing about New Cuyama because it combines the count of people, the elevation (in feet above sea level), and the year the town was established (on a certain calendar). Having three numbers does not ensure that addition will produce any sensible result.

What Is Missing from Stevens

Stevens' four levels are usually presented in the geographic literature as a complete set, but they are not enough for practical applications of geographic information. There are at least four revisions that need to be added.

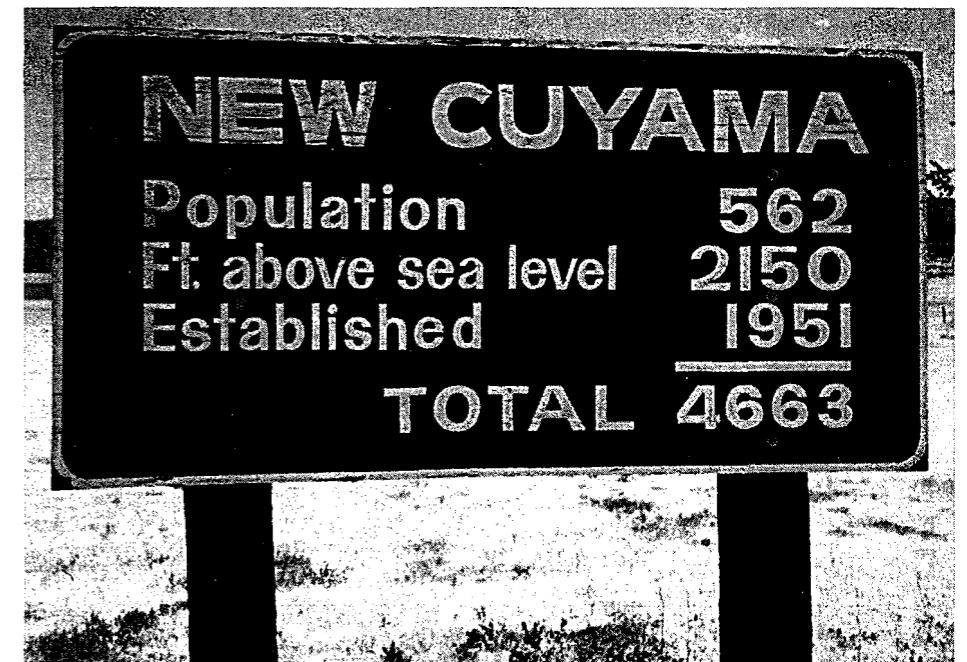


Figure 1-10 Sign posted at the entry to New Cuyama, California.

Absolute Scales Ratio is not actually the highest level of measurement. A ratio scale is higher than an interval scale because the value of zero is no longer arbitrary. A higher level of measurement would be achieved if the unit of measure were not arbitrary. When the whole scale is predetermined or *absolute* (Ellis 1966), no transformations that preserve the meaning of the measurement can be made. One example of such an absolute scale is *probability*, where the meaning of zero and one are given. *Even though it is common to report probabilities as percentages, the relationships of probability (such as Bayes' law of conditional probability) operate correctly only when scaled from zero to one.*

Cyclical Measures While Stevens' levels deal with an unbounded number line, there are many measures that are bounded within a range and repeat in some cyclical manner. Angles seem to be ratio, in the sense that there is a zero and an arbitrary unit of measure (degrees, grads, or radians); however, angles return to their origin. The direction 359° is as far from 0° as 1° is. Any general measurement scheme needs to recognize the existence of such cyclical measures.

Counts Another class of geographic measurements deals with counting objects aggregated over some region in space, such as a human population. The objects counted are discrete; there is no half person. Yet, unlike the other discrete levels of measurement (nominal and ordinal), the result of a count is a number. Since the zero is a fixed value, counts may seem to be ratios, but the units of a count are not arbitrary, so they cannot be rescaled as freely. Counts are more similar to absolute scales, with a restriction to discrete integers. They become ratios when the unit of measure is rescaled to "population in millions" or something that loses the discrete identity of the objects.

Graded Membership in Categories A further criticism of Stevens' system is that nominal categories are not always as simple as portrayed. Nominal measures apply the strict rules of classical set theory. All members of a set are meant to belong in that category equally; these sets are called "sharp." Many classifications, however, adopt more flexible rules; they involve some kind of graded memberships as formalized in **fuzzy set theory** or they involve comparison to a **prototype** member of the class. In both these situations, an object will have some degree of membership (represented by a proportion or percentage), rather than just belonging or not. Some members of the group are just more typical than others. Accommodating a more nuanced

Fuzzy set theory: An extension to set theory that permits an object to have a degree of membership (usually represented as a number between 0 and 1). Fuzzy membership values do not have to follow the rules of probability.

Prototype: An approach to categorization that defines a category by identifying a particular object as the typical example. Other objects assigned to this category may not share all characteristics with the prototype object. The degree of resemblance represents graded membership.

approach to categories remains a research frontier in GIS (Burrough and Frank 1996).

Thus, Stevens' four levels of measurement are not the end of the story. A closed list of levels arranged on a progression from simple to more complex does not cover the diversity of geographic measurement. Still, Stevens' terminology provides a starting point for the bulk of common situations.

Applying Levels of Measurement to Attribute Reference Systems

The previous section ended with the conclusion that attribute reference systems seemed too varied for standardized treatment. Though they do not specify all details, Stevens' four levels of measurement (with extensions) do prescribe the information required for an attribute reference system (Table 1-3). Absolute measurements can simply state what they measure because the whole scheme is implicit. For a count, the reference system is the kind of object counted. For a ratio scale, the unit of measure must be given, along with some additional information to sort out the subcases of cyclical scales and derived ratios. For interval measures, the units and the zero point are required. The categorical levels require more information because each category has its own definition. Perhaps the ordinal categories of "somewhat poorly drained" and "poorly drained," for example, are divided at a specific threshold on a ratio measure of permeability. Other ordinal values may not have explicit links to a numerical scale, just a ranking. With nominal categories, each category needs to be described. Some category systems are simply lists, as in the Anderson land use codes (Anderson et al. 1976). Another way to present categories uses a series of structured questions. For example, does the tree have leaves or needles? Are the needles in groups or singly attached? Are the groups of five, three, or two? Such a key can emphasize different characteristics in the various paths, each leading to a particular category.

TABLE 1-3 Information Content for Attribute Reference Systems

<u>Level of Measurement</u>	<u>Information Required</u>
Nominal	Definitions of categories
Graded membership	Definition of categories plus degree of membership or distance from prototype
Ordinal	Definitions of categories plus ordering
Interval	Unit of measure plus zero point
Extensive ratio	Unit of measure (additive rule applies)
Cyclic ratio	Unit of measure plus length of cycle
Derived ratio	Unit of measure (ratio of units; weighting rule)
Counts	Definition of objects counted
Absolute	Type (probability, proportion, etc.)

Attribute Reference Systems for the La Selva Project

The La Selva project demonstrates a diversity of attribute reference systems. As in many GIS projects, the bulk of the data sources available to the students for the La Selva project were in categories. The project focused on the Sarapiquí Annex, a parcel of land. Inside this boundary, the map showed some points representing the pipes that demarcate the spatial reference system and other points depicting stumps (signs of logging activity). A tree stump is a member of a very simple nominal category. There were also roads, trails, and streams. To some extent, the roads and trails are a part of an ordinal set of classes. The primary content of the map was a land use delineation in which the categories were ordered along a gradient of human disturbance. "Primary Forest" included the rain forest with the least human influence, followed by "Selectively Cut," "Cleared Land," and so on. The land-use mappers did not document the rules that were applied in deciding how many trees had to be removed to qualify for each category. Presumably, these categories made sense considering the local economy. Some forest was completely cleared to create agriculture, while other operations targeted specific kinds of trees.

The La Selva project also used some information from the Reserve's existing database. Elevation data appeared in standard meters. The zero was mean sea level, which was not particularly relevant in the cloud forest. For all practical purposes, the elevation data were interval measures, largely useful when compared to each other to construct measures (to be introduced in Chapter 7) such as slope gradient (an absolute scale) or slope aspect (a cyclical ratio). The soils data were classified according to an international nomenclature for soils. These classes were divided from each other according to multidimensional thresholds: permeability, organic content, and grain sizes. Analytically they were treated as sharply distinct sets, although they probably represent a series of complex gradients.

SUMMARY

Geographic information must be embedded in a reference system for time, space, and attribute. Time and space have fairly standardized systems in common use. Attributes, by contrast, come in all flavors. The generic typology of Stevens' levels of measurement provide a starting point to develop the information content required. Numerous additions and special cases must be recognized.

To use measurements effectively, additional distinctions must be made. These distinctions do not come from the numbers but from a larger framework surround-

ing the measurements. A spatial reference system provides a mechanism to construct a more integrated structure, but coordinates by themselves do not ensure compatibility of diverse information. Time and attribute also have reference systems, but these three systems just provide the basic axes. A more comprehensive framework must include the interactions of these three components. The next chapter develops such a framework for geographic information.

REPRESENTATION

CHAPTER OVERVIEW

- Introduce the primitives used to represent spatial and attribute measurements.
- Describe basic representation models: vector and raster.
- Follow steps that convert existing documents into digital databases.
- Introduce components of data quality evaluation applied to digitizing.

Representation involves a symbol acting in the place of some entity. The clearest precedent for geographic representation involves the graphic symbols used on traditional maps. A map populates a small space with the representation of a larger space, using map symbols to stand in place of things in the world. Thus, a small star serves to encode that a particular city is the capital of its state. While the direct representation of space facilitates visual interpretation, the compression to a smaller space creates many dilemmas. Traditional map symbolism often dictates a measurement framework to suit the technical limits of graphic reproduction. This legacy still influences the understanding of geographic information.

Digital representations are far less direct than maps in structuring spatial relationships. Digital symbols are extremely simple until built into more complex structures. Despite their simplicity, digital representations improve on the amazing capabilities of the printed map. This chapter reviews the primitive elements used to represent geographic data in digital media and then describes the basic data models and **data structures** that organize these primitives into a useful representation. The final portion of the chapter traces the various steps used to convert data sources into an information system.

Data structure: Arrangement of data entities that permits the construction of relationships through software operations; implements a data model.

PRIMITIVES FOR REPRESENTATION

Computers provide a few methods to represent numbers built into the hardware. The basic unit of storage is a bit, a single binary digit. These bits are grouped into larger units to represent numbers. Most computers provide three “word” formats for numbers: *integer*, *floating point*, and *double precision* (Figure 3-1). At one time, the number of bits in a word varied, but these days most computers store integer and floating-point values in 32 bits and double precision in 64 bits. Thirty-two-bit integers can represent any whole number from -2147483647 to 2147483648 . With integers, results of any division that involves a remainder will be truncated to the next lowest integer (by absolute value). Floating point uses a logarithmic notation with an exponent to establish magnitude and a mantissa for the value. This format provides about six decimal digits of resolution, scaled over a much larger range (typically -10^{36} to 10^{36}). Calculations in floating point are more forgiving than the truncation associated with integers, but ultimately a floating-point number is still represented with finite resolution. Double precision provides about 15 digits, essentially an extended form of floating point, but storage is doubled and speed of calculations is usually decreased. Choosing between these formats has consequences in handling spatial information.

Besides numbers, computer representation can manage information coded as characters. Most storage media work in 8-bit *bytes*. International standards provide a coding for certain basic characters, though the full range of special characters remains far from standardized. Since characters are nominal symbols, the hardware

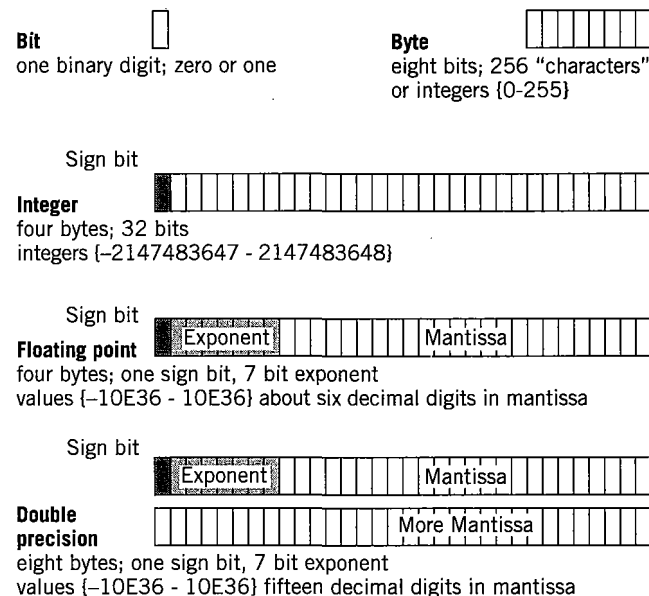


Figure 3-1 Computer storage at the most basic level.

can do little other than copy them. The primitive of a byte handles a single character. Groups of characters require rudimentary data structures, for example, a common way to store a *string* of characters starts with a byte count followed by that many bytes.

Primitives for Attributes

Attribute values in a digital system are directly encoded using the units of computer storage. The case and variable framework of the geographical matrix encourages the data structure of a two-dimensional array. One storage location is allocated for each attribute value. Numbers are often encoded in floating point because six digits are usually adequate for the significant portion of an attribute measurement. Counts are usually encoded as integers to make use of the increased range, although a 32-bit integer is no longer adequate to represent the sum of the world population of humans. When storage bulk becomes a concern, as with remotely sensed imagery, continuous measurements may be rescaled from their original range into a range from 0 to 255 and encoded as an integer in a single byte. Using this representation, an attribute that was conceptually continuous has become effectively discrete. Most users forget these details and simply treat an attribute as a “number”; they assume that the computer can handle all the intuitively obvious arithmetic operations. Computer hardware manipulates the storage units without recognizing different levels of measurement; thus, the validity of calculations depends on external meaning.

Nominal attributes must be coded as some kind of number as well. The coding rules become a part of the data structure schema that must be preserved to retain the meaning of the information. Categories may be encoded in a single byte with no loss of information content, as long as there are only 256 possibilities.

Primitives for Time

Time can be treated simply as another attribute for most purposes. Representation of dates became an issue of great importance around the year 2000. The Y2K problem arose from a short-sighted choice to represent dates with only two digits for years. This was a very common simplification in accounting and other business applications. Large amounts were spent trying to revise programs so that a difference in dates would not come up with erroneous negative results. For example, $2001 - 1980 = 21$ years, but $01 - 80 = -79$. Unless carefully designed, choices of representation create unintended consequences.

Primitives for Space: Coordinates

Representation of space requires a spatial reference system as defined in Chapter 1. On a map, location is encoded by relative position on the piece of paper. In a computer database, location is based on analytical geometry. For practical purposes, the spatial component of geographic information is represented in the form of coordi-

nates, that is, ordered measurements relative to a spatial reference system. These measurements may be angles on an ellipsoid (latitude, longitude) or orthogonal distances on a projection plane (Figure 1-3). There are different ways to encode equivalent measurements, even on a plane. For example, a point can be recorded either by two distances along X and Y axes or as one distance and one angle, using a radial coordinate notation (Figure 3-2). Radial coordinates preserve all the information content, but organize it differently. Each GIS software package will implement a certain range of alternative spatial reference systems with associated representations. These possibilities are strongly governed by conventions. Cartesian coordinates are usually stored in the order X then Y with positive values extending “right” and “up,” but many raster representations (introduced below) have the origin at the top of the screen and count lines “down” instead. There is even more variability in representing latitude and longitude. Most follow the convention that places latitude before longitude, but others associate longitude with the X axis and thus place it first. The old Babylonian base-60 remains in use for degrees, minutes, and seconds, but some software uses decimal degrees (decimal fractions of degrees) while others use radians (and decimal fractions). All of these representations have simple mathematical relationships due to different units of measure, but sometimes a display will seem to be in the wrong scale or flipped in mirror image until the differences in conventions are smoothed out.

Coordinates are almost invariably stored as pairs of the basic computer words, with interesting consequences (Chrisman 1984a). All the number primitives supported in computer hardware (integer, floating point, and double precision) have limited resolution because they are effectively integers (Figure 3-1). Limited resolution has additional consequences when applied to geometric representation; integers cannot represent all intermediate locations along diagonal lines (Figure 3-3). Most diagonal lines from one whole number pair to another pass through locations that cannot be represented in the integer system without bending the straight-line segment (Egenhofer and Herring 1991). In order to limit these geometric problems, the representation systems for coordinates should store more resolution than the application seems to require. Integer storage also means that no coordinate representation can offer truly continuous variation. The size of the discrete jumps in resolution can be made so tiny that they seem effectively continuous, but only at some expense.

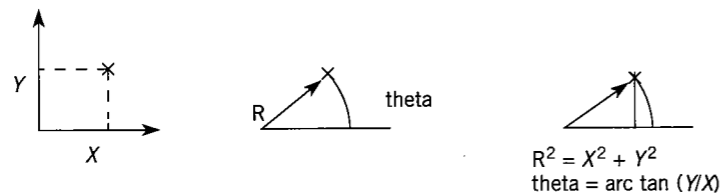


Figure 3-2 Conversion between measurements in two dimensions. Radial coordinates (angle and distance from an origin) are equivalent to cartesian coordinates (distances parallel to two orthogonal axes). The angle measurement is cyclical, not extensive.

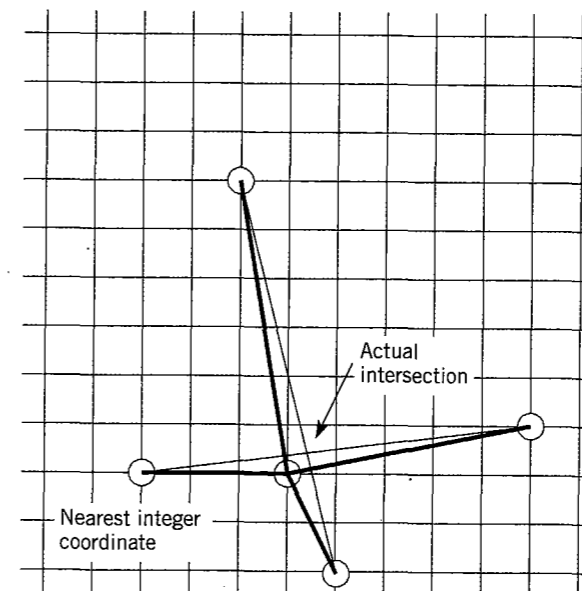


Figure 3-3 Representing intermediate locations in an integer coordinate system. Thicker gray lines represent the result of intersection of two straight lines once rounded off into the integer space. Note that the new segments have different slopes from their parent segments. Coordinates stored as floating-point numbers do not avoid this problem.

Resolution has a consequence in setting the maximum extent of the area represented. Single-word floating point with its six digits can circle the equator (just over 40,000 km) so that kilometers can be resolved, but not much more. On a county scale, these words may resolve centimeters, but only by placing a shifted false origin near the county (also termed a *local offset*). Otherwise, storage bits are wasted in measuring each coordinate relative to some distant origin. With a local offset, floating point can be used for most county-sized databases. Unfortunately, however, the database of the adjacent county might not be appended without overflowing the range of a floating-point value. Some software has opted for double precision, paying the penalty of slower speed and increased storage to avoid the troubles of single-word storage. With double precision, there is enough resolution to distinguish the left and right side of every virus particle on the planet. For mapping purposes, this resolution is spurious and must be filtered back to some reasonable level. It is particularly odd to use such profligate storage for measurements obtained from inaccurate maps by fallible digitizing systems.

REPRESENTATION MODELS AND DATA STRUCTURES

Representation of geographic information proceeds by organizing the primitives into more complex structures. These data structures often provide the key technical differences between competing software packages. The specific details of a data structure act as instances of a more generic data model of entities and their relationships. Much of the substance of a data model comes from a measurement framework, as

discussed in the previous chapter. Many alternative data structures are possible, but there are relatively few generic models behind them. This section will concentrate on the two dominant models of geometric representation in GIS, *vector* and *raster*, followed by consideration of database architecture for the whole system.

Vector Model

Measurement frameworks based on attribute control are implemented most directly using the **vector** model. Based on analytical geometry, a vector model builds a complex representation from primitive objects for the spatial dimensions: points, lines, and areas. These primitives have a nested dependency: areas are described by boundary lines, and the location for a line can be approximated by a string of line segments connecting a series of points. At the base, points are represented by coordinates. Cartographic data structures usually do not provide more complicated options between points. This simplicity contrasts with the richness of different curves in drafting or illustration software, where the paths between points may be **Bézier curves** or **splines**, not just straight-line segments. Of course, these software packages are oriented toward display, not analytical operations using the data. In engineering practice, many highway features are laid out with circular arcs or conic spirals, but most natural features do not have a preferred mathematical curve. Software systems that allow different curves to connect points have to provide much more complex processing for relationships such as geometric intersections. For general cartographic representation, segments provide a simple, versatile approximation.

Representing Isolated Objects The spatial object measurement framework translates into a simple vector representation with each line and polygon defined by a string of coordinates (Figure 3-4). Any string can represent a line, but a polygon should close. Collections of these spatial objects are stored in what are often termed *shape files*. As long as the objects remain isolated, which is an axiom of the measurement framework, the representation serves its basic purpose. To represent an inner ring inside a polygon, the isolated data structure often adopts a convention of a *retraced line* connecting the outer ring to the inner ring. Instead of having two disjoint rings, at same point around the outer ring, a line is inserted connecting to the inner ring. Once the inner ring has been traced, the inserted line is repeated exactly in the reverse direction (thus it is retraced). Software using this representation has to check for these retraced segments so that they are not drawn on the display. For example, in Figure 3-4, a shape describing South Africa would require a retraced line

Vector: A spatial data model based on geometric primitives (point, line, and area), located by coordinate measurements in a spatial reference system; from mathematical term for a direction, or a directed line segment.

Bézier curve: A smooth curve that passes through specified points with a given direction (tangent) at those points.

Spline: A smooth curve that models the behavior of a thin spring (with a given modulus of elasticity) constrained to pass through specified points.

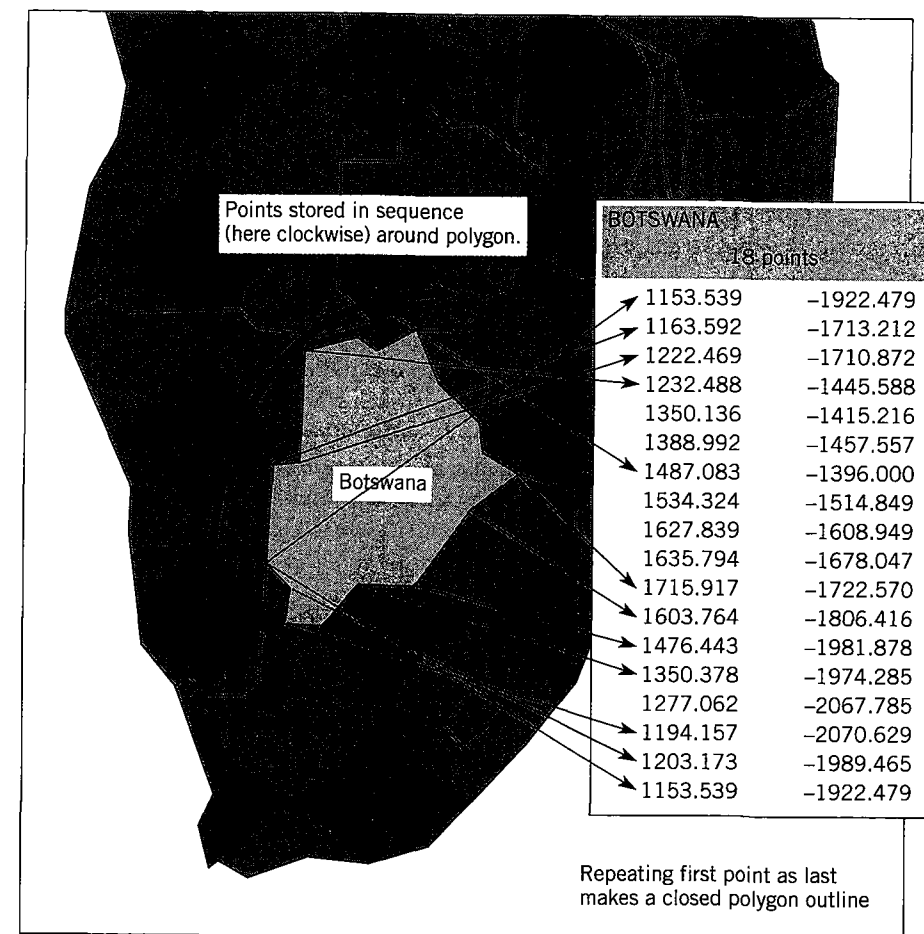


Figure 3-4 Spatial object framework implemented in a vector representation. Each object contains its own list of points (coordinates). No relationships are stored.

to remove Lesotho from what would otherwise be inside the outer ring of South Africa. What appears to be a simple data structure can cause unexpected software complexity.

The isoline framework, since it uses isolated contour lines, also translates directly to this simple vector representation. In addition to the closed contour loops, an isoline data structure might decide to represent the relationships between adjacent (nested) lines up and downhill.

There is considerable potential for inconsistency if an isolated object representation is applied to a connected network of polygons. Since each object is independently encoded as a closed loop, each boundary is represented twice. For example, the first segment of the Botswana polygon (in Figure 3-4) would also appear somewhere in the shape for Namibia. A gap or overlap could easily occur without any easy

method to detect it (Figure 3-5). Software built upon shape files often includes topological processing to avoid such inconsistencies.

Representing Topological Relationships The topological data model has been more commonly used in software that implements a full range of operations on vector representations. The topological model incorporates network relationships along with the coordinate measurements (Figure 3-6); thus, it can handle the requirements of the connected coverage frameworks. This model centers around the boundary with explicit connection to nodes at each end as well as to polygons on the left and right (as introduced in Chapter 2). These relationships can be implemented in a number of different specific data structures, particularly in relating polygons to their boundaries and inner rings to their outer rings (Gold 1988). One common structure creates a variable-length list of the chains around each ring, with notation for direction, as shown in Figure 3-6. Another approach, called winged-edge, stores the “next clockwise” chain pointer at the ends of chains. With a single starting point, a ring can be assembled by following from pointer to pointer. However they are implemented, the data structures provide access to the same relationships.

The important characteristic of all the vector methods is that they permit essentially free placement of point locations and of boundary lines to represent categories. Thus, the vector model directly implements the intent of attribute-controlled measurement frameworks. Because of this linkage, the measurement framework and the representation may seem inseparable, but representation remains a distinct choice. The vector model can be applied to other measurement frameworks, such as the tri-

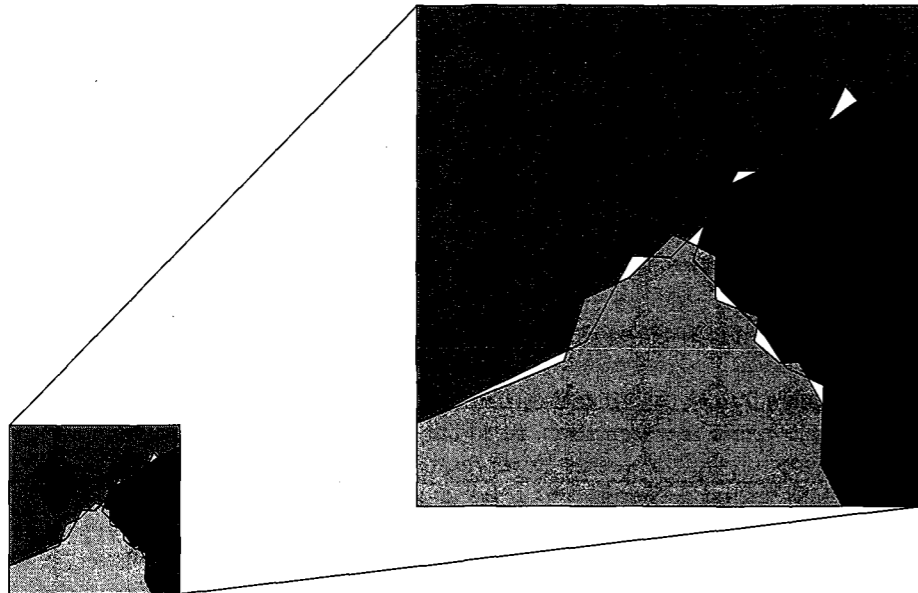


Figure 3-5 Adjacent polygons represented by isolated boundary lines may create slivers (gaps and overlaps) if adjacent boundaries are not identical [much enlarged].

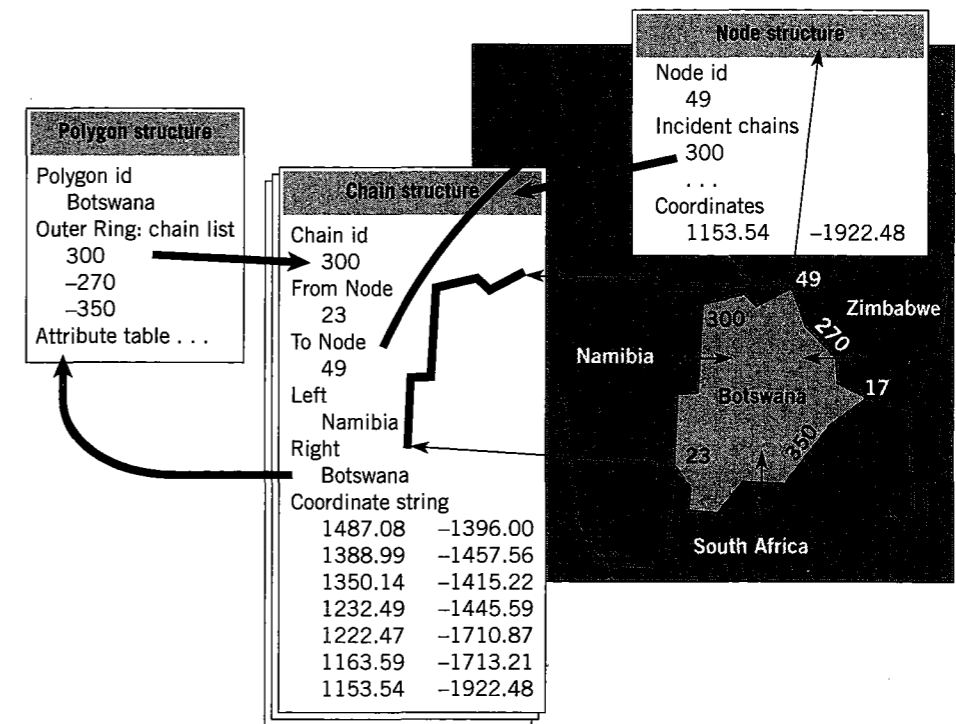


Figure 3-6 Topological data structure includes relationships between the components of a connected network. The polygon representing Botswana has three neighbors (Namibia, Zimbabwe, and South Africa) at this scale. (Actually, the node on the Zambesi River includes a 100-meter border with Zambia, not representable at this scale.) Each neighbor requires a border to separate the two countries. Each border chain begins and ends at a node. This diagram shows a polygon with a list of chains around its outer ring (using a sign for reverse direction).

angles of a TIN. Vector data structures are the method of choice for choropleth mapping of continuous attributes. As long as geographic data entry is difficult and expensive, there is a strong incentive to use one geometric description for many attributes, the approach of the geographical matrix. The evaluation of competing approaches will be revisited in Chapter 10.

Raster Model

The other major family of representation models is called **raster**. Whereas the vector model is constructed from geometric primitives as a logical structure, the raster model has close links to the physical characteristics of computer graphics hardware. Raster derives from a word used in mechanical engineering for a tool that sweeps

Raster: A spatial data model based upon a regular tessellation of a surface into pixels or grid cells.

back and forth as it advances. Television technology uses an electronic sweep gun, so the term came to refer to the rows on the screen and, eventually, to the cellular nature of a cathode ray tube (CRT) display. The hardware structure of remote-sensing sensors, of line printers, and of CRT displays all contributed to the popularity of the raster structure. In addition to the hardware connection, raster structures could be implemented easily with the elementary data structures available even in the earliest programming languages.

The raster model divides the region into rectangular building blocks (grid cells or *pixels*) that are filled with the measured attribute values. The raster approach is directly related to the frameworks that control space in order to measure attributes. Raster cells are located within a spatial reference system, but they deliberately limit resolution to act as control. In many cases, the raster geometry is specified by the original sensing hardware (Figure 3-7). Raster representation is not restricted to its related measurement framework; for example, attribute-controlled measurements can be represented in a raster; in effect the process creates a composite using both attribute and spatial control successively.

An array in computer storage (random access memory or disk) provides the most direct implementation of a raster representation. At sufficiently high resolution, a raster representation will have many adjacent cells with identical values, particularly when applied to categorical attributes. There are many methods of **compression** possible. Overall, compression can be divided into *loss-less* methods, which preserve all the data, and *statistical* methods, which might change some data to simplify the encoding. The latter group effectively uses a form of cartographic generalization (usually implemented using the neighborhood operations described in Chapter 7), followed by one of the loss-less methods.

There are a number of loss-less compression methods; a few are diagrammed in Figure 3-8. If a raster contains chunks of cells with identical values, *run length encoding* provides significant compression. Instead of storing each cell, each component stores a value and a count of cells along the row sharing that value (Figure 3-8a). If there is only one cell, the storage doubles (from one byte to two), but for three or more there is a reduction. In the special case of a binary (black/white) image, the compression need only store the cells where the value changes. Run length methods work along one row. More advanced methods store differences between rows instead of treating each row separately (Figure 3-8b). Compression algorithms form a routine part of the **TIFF** standard used by facsimile machines to reduce telecommunications costs, but methods that work for a digital transmission may not be best tuned for geographic access. Various forms of **quadtree** data structures try to take advan-

Compression: A software procedure that encodes a data structure so that its storage occupies less space (under certain conditions); may preserve all the information (loss-less) or deliberately simplify.

TIFF: Tagged Image File Format: a family of image encoding formats that can vary the resolution and the number of bits used to represent each cell.

Quadtree: A spatial data structure that organizes a hierarchical structure of square cells through iterative division into four daughter cells (Samet 1990).

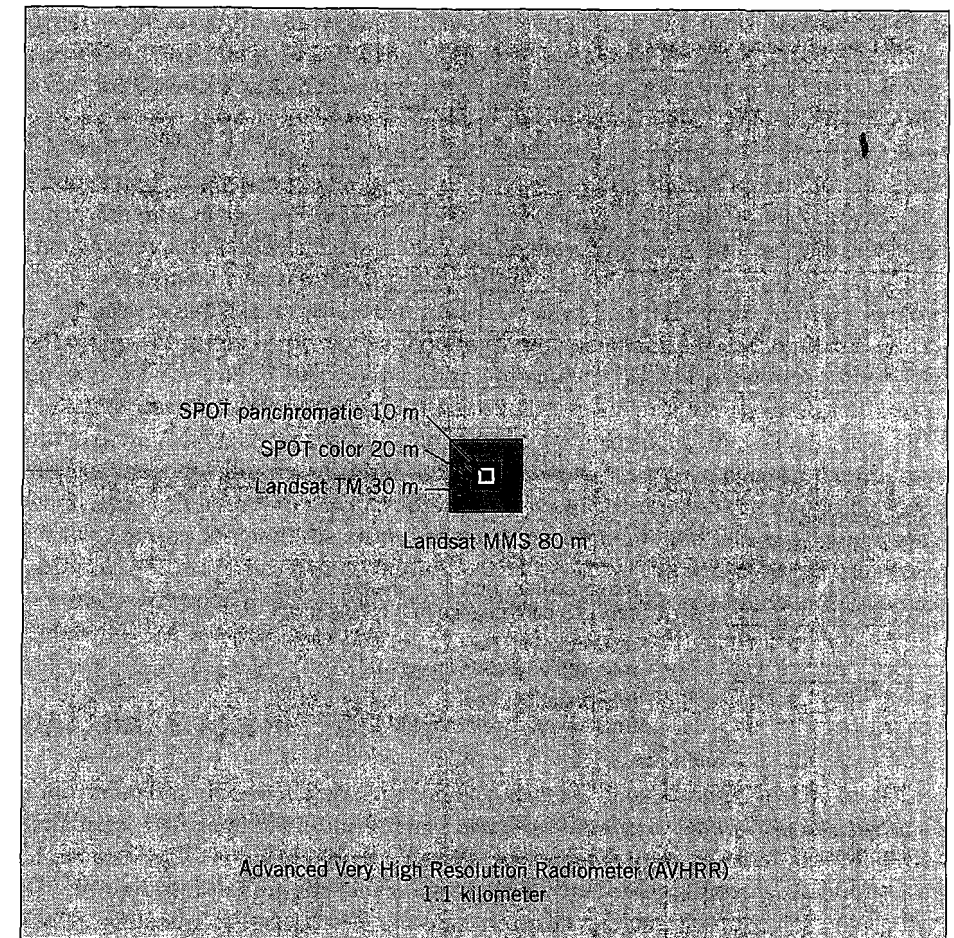


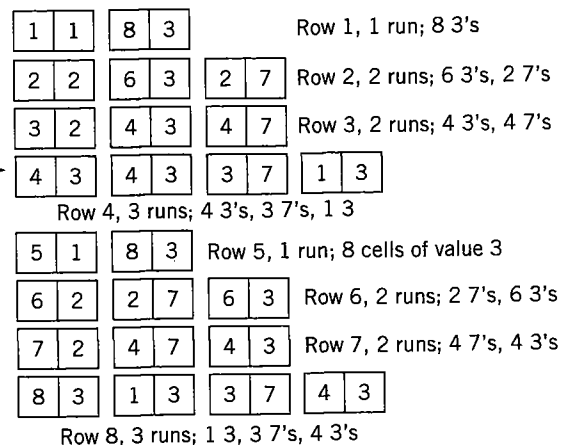
Figure 3-7 Pixel sizes from common satellite sensing systems.

tage of two-dimensional character of spatial data for compression. A quadtree works by iterative division of a region into four square subunits. A *region quadtree* stores a categorical attribute and will not subdivide a square that is homogeneous (Figure 3-8c). The pattern of this compression also provides a measure of spatial variability (Csillag and Kummert 1990).

In addition to compression, a key issue with a raster system is the size of the cells. Smaller cells permit the raster to approximate the flexibility of the vector system as closely as required, but at a price in the storage consumed. A coarse cell system is sometimes distinguished from raster and may be called a *grid cell* system. A grid cell is sufficiently large that it can no longer be treated as a point, so one of the area-based rules discussed in Chapter 2 must be used to make a measurement. Though all raster pixels actually occupy some space, they should be so small that there is no internal

(a) Compression by Run Length Encoding (along rows)

3	3	3	3	3	3	3	3
3	3	3	3	3	3	7	7
3	3	3	3	7	7	7	7
3	3	3	3	7	7	7	3
3	3	3	3	3	3	3	3
7	7	3	3	3	3	3	3
7	7	7	7	3	3	3	3
3	7	7	7	3	3	3	3



(b) Compression by Row Differences

1	1	1	8	3
2	1	7	2	7
3	1	5	2	7

Interpretation: Row 1, 1 change section; from cell 1, 8 cells change to 3 (whole row)

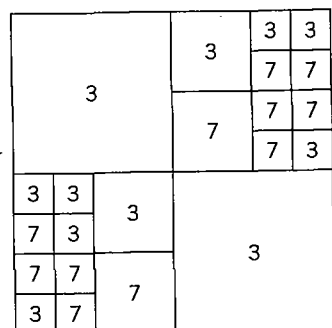
Interpretation: Row 2, 1 change section; from cell 7, 2 cells change to 7

Interpretation: Row 3, 1 change section; from cell 5, 2 cells change to 7

and so on...

(c) Compression by Quadrees

3	3	3	3	3	3	3	3
3	3	3	3	3	3	7	7
3	3	3	3	7	7	7	7
3	3	3	3	7	7	7	3
3	3	3	3	3	3	3	3
7	7	3	3	3	3	3	3
7	7	7	7	3	3	3	3
3	7	7	7	3	3	3	3



64 cells compress to 22 quadtree leaves.

Figure 3-8 Compression of raster representations. The repeated values in adjacent cells can be reduced along the rows, between the rows, or by quadtree.

detail at the resolution of representation. For tiny pixels, the measurement rule may not matter very much. Many of the original systems used for natural resource inventories in the 1970s used quite crude cell sizes. For example, the Land Use and Natural Resource project in New York State used 1-km squares; the Maryland Automated Geographic Inventory used 2000-foot cells. Crude cell sizes may have disappeared from state planning agencies, but they remain quite common in global environmental models. Cell sizes of 1° by 1° or even 5° by 5° are quite common, even though they are hardly square or uniform as they go poleward. For cells this large, the measurement rule becomes quite critical.

In the late 1970s, a serious and heated debate occurred over the virtues and disadvantages of raster versus vector models of representation. At that time, the debate often revolved around efficiency in implementation of data structures rather than contrasting the fundamental models. The choices made in the early days of GIS were often driven by technology, not application. In retrospect, the debate made the choices seem needlessly exclusive; each representation serves a measurement framework, and each framework has its appropriate uses. It is much more important to tailor the representation to the axioms of the measurement (and ultimately to the purposes of the enterprise) than to let the technology drive the decisions. As computing power has dropped in price, software has been able to deliver a much closer approximation of the conceptual models.

Database schema: Logical arrangement of tables, attributes, and integrity rules to structure a database. Involves definitions of entities and their relationships.