

3 Geodesy, Datums, Map Projections, and Coordinate Systems

Introduction

Geographic information systems are different from other information systems because they contain spatial data. These spatial data include coordinates that define the location, shape, and extent of geographic objects. To effectively use GIS, we must develop a clear understanding of how coordinate systems are established for the Earth, how these coordinates are measured on the Earth's curving surface, and how these coordinates are transferred to flat maps. This chapter introduces *geodesy*, the science of measuring the shape of the Earth, and *map projections*, the transformation of coordinate locations from the Earth's curved surface onto flat maps.

Defining coordinates for the Earth's surface is complicated by three main factors. First, most people best understand geography in a Cartesian coordinate system on a flat surface. Humans naturally perceive the Earth's surface as flat, because at human scales the Earth's curvature is barely perceptible. Humans have been using flat maps for more than 40 centuries, and although globes are quite useful for perception and visualization at extremely small scales, they are not practical for most purposes.

A flat map must distort geometry in some way because the Earth is curved. When we plot latitude and longitude coordinates on a Cartesian system, "straight" lines will appear bent, and polygons will be

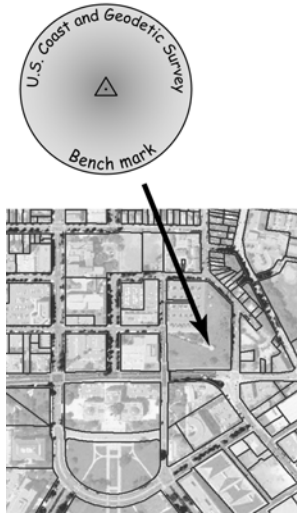
distorted. This distortion may be difficult to detect on detailed maps that cover a small area, but the distortion is quite apparent on large-area maps. Because measurements on maps are affected by the distortion, we must somehow reconcile the portrayal of the Earth's truly curved surface onto a flat surface.

The second main problem in defining a coordinate system results from the irregular shape of the Earth. We learn early on that the Earth is shaped as a sphere. This is a valid approximation for many uses, however, it is only an approximation. Past and present natural forces yield an irregularly shaped Earth. These deformations affect how we best map the surface of the Earth, and how we define Cartesian coordinate systems for mapping and GIS.

Thirdly, our measurements are rarely perfect, and this applies when measuring both the shape of the Earth, and the exact position of features on it. All locations depend on measurements that contain some error, and on analyses that must make some assumptions. Our measurements improve through time, and so does the sophistication of our models, so our positional estimates improve; this evolution means our estimates of positions change through time.

Because of these three factors, we often have several different sets of coordinates to define the same location on the surface of the Earth. Remember, coordinates

Coordinates for a Point Location



From Surveyor Data:

	Latitude (N)	Longitude (W)
NAD83(2007)	44 57 23.23074	093 05 58.28007
NAD83(1986)	44 57 23.22405	093 05 58.27471
NAD83(1996)	44 57 23.23047	093 05 58.27944

	X	Y	
SPC MNS	317,778.887	871,048.844	MT
SPC MNS	1,042,579.57	2,857,766.08	sFT
UTM15	4,978,117.714	492,150.186	MT

From Data Layers:

	X	Y	
MN-Ramsey	573,475.592	160,414.122	sFT
MN-Ramsey	174,195.315	48,893.966	MT
SPC MNC	890,795.838	95,819.779	MT
SPC MNC	2,922,552.206	314,365.207	sFT
LCC	542,153.586	18,266.334	MT

Figure 3-1: An example of different coordinate values for the same point. We may look up the coordinates for a well-surveyed point, and we may also obtain the coordinates for the same point from a number of different data layers. We often find multiple latitude/longitude values (surveyor data, top), or x and y values for the same point (surveyor data, or from data layers, bottom).

are sets of numbers that unambiguously define locations. They are usually x and y values, or perhaps x, y, and z values, or latitude and longitude values unique to a location. But these values are only “unique” to the location for a specified set of measurements and time. The coordinates depend on how we translate points from a curved Earth to a flat map surface (first factor, above), the estimate we use for the real shape of the Earth (second factor), and what set of measurements we reference our coordinates to (the third factor). We may, and often do, address these three factors in a number of different ways, and the coordinates for the same point will be different for these different choices.

An example will help clarify this concept. Figure 3-1 shows the location of a U.S. bench mark, a precisely surveyed and monumented point. Coordinates for this point are maintained by Federal and State government surveyors, and resulting coordinates shown at the top right of the figure. Note that there are three different versions of the latitude/

longitude location for this point. In this case, the three versions differ primarily due to differences in the measurements used to establish the point’s location, and how measurement errors were adjusted (the third factor, discussed above). The GIS practitioner may well ask, which latitude/longitude pair should I use? This chapter contains the information that should allow you to choose wisely.

Note that there are also several versions of the x and y coordinates for the point in Figure 3-1. The difference in the coordinate values are too great to be due solely to measurement errors. They are due primarily to how we choose to project from the curved Earth to a flat map (the first factor), and in part to the Earth shape we adopt and the measurement system we use (the second and third factors).

We first must define a specific coordinate system, meaning we choose a specific way to address the three main factors of projection distortion, an irregularly shaped Earth, and measurement imprecision. There-

after the coordinates for a given point are fixed, as are the spatial relationships to other measured points. But it is crucial to realize that different ways of addressing 1) the Earth's curvature, 2) the Earth's deviation from our idealized shape, and 3) inevitable inaccuracies in measurement, will result in different coordinate systems, and these differences are the root of much confusion and many errors in spatial analysis. As a rule, you should understand the coordinate system used for all of your data, and convert all data to the same coordinate system prior to analysis. The remainder of this chapter describes how we define, measure, and convert among coordinate systems.

Early Measurements

In specifying a coordinate system, we must first define the size and shape of the Earth. Humans have long speculated on this. Babylonians believed the Earth was a flat disk floating in an endless ocean, a notion adopted by Homer, one of the more widely known Greek writers. The Greeks were early champions of geometry, and they had many competing views of the shape of the Earth. One early Greek, Anaximenes, believed the Earth was a rectangular box, while Pythagoras and later Aristotle reasoned that the Earth must be a sphere. He observed that ships disappeared over the horizon, the moon

appeared to be a sphere, that the stars moved in circular patterns, and that constellations shift when viewed from different ends of the Mediterranean Sea. These observations were all consistent with a spherical Earth.

The Greeks next turned toward estimating the size of the sphere. The early Greeks measured locations on the Earth's surface relative to the Sun or stars, reasoning they provided a stable reference frame. This assumption underlies most geodetic observations taken over the past 2000 years, and still applies today, with suitable refinements.

Eratosthenes, a Greek scholar in Egypt, performed one of the earliest well-founded measurements of the Earth's circumference. He noticed that on the summer solstice the Sun at noon shone to the bottom of a deep well in Syene. He believed that the well was located on the Tropic of Cancer, so that the Sun would be exactly overhead during the summer solstice. He also observed that 805 km north in Alexandria, at exactly the same date and time, a vertical post cast a shadow. The shadow/post combination defined an angle which was about $7^{\circ}12'$, or about 1/50th of a circle (Figure 3-2).

Eratosthenes deduced that the Earth must be 805 multiplied by 50, or about 40,250 kilometers in circumference. His calculations were all in stadia, the unit of measure of the time, and have been converted

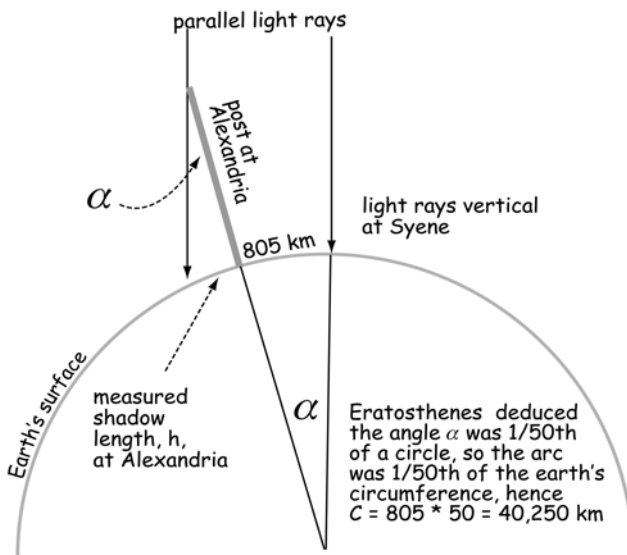
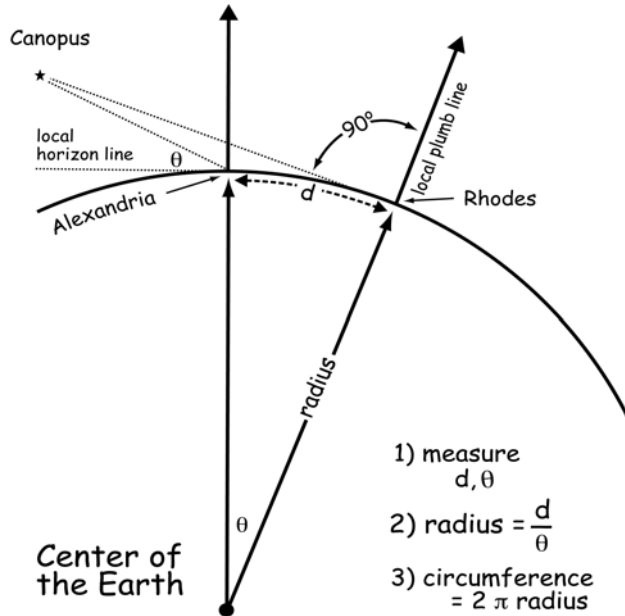


Figure 3-2: Measurements made by Eratosthenes to determine the circumference of the Earth.

Figure 3-3: Posidonius approximated the Earth's radius by simultaneous measurement of zenith angles at two points. Two points are separated by an arc distance d measured on the Earth surface. These points also span an angle θ defined at the Earth center. The Earth radius is related to d and θ . Once the radius is calculated, the Earth circumference may be determined. Note this is an approximation, not an exact estimate, but was appropriate for the measurements available at the time (adapted from Smith, 1997).



here to the metric equivalent, using our best idea of the length of a stadia. Eratosthenes' estimate differs from our modern measurements of the Earth's circumference by less than 4%.

The accuracy of Eratosthenes' estimate is quite remarkable, given the equipment for measuring distance and angles at that time, and because a number of his assumptions were incorrect. The well at Syene was located about 60 kilometers off the Tropic of Cancer, so the Sun was not directly overhead. The true distance between the well location and Alexandria was about 729 kilometers, not 805, and the well was 3°3' east of the meridian of Alexandria, and not due north. However these errors either compensated for or were offset by measurement errors to end up with an amazingly accurate estimate.

Posidonius, another Greek scholar, made an independent estimate of the size of the Earth by measuring angles from local vertical (plumb) lines to a star near the horizon (Figure 3-3). Stars visible in the night sky define a uniform reference. The angle between a plumb line and a star location is called a *zenith angle*. The zenith angle can be measured simultaneously at two locations

on Earth, and the difference between the two zenith angles can be used to calculate the circumference of the Earth. Figure 3-3 illustrates the observation by Posidonius at Rhodes. The star named Canopus was on the horizon at Rhodes, meaning the zenith angle at Rhodes was 90 degrees. He also noticed Canopus was above the horizon at Alexandria, meaning the zenith angle was less than 90 degrees. The surface distance between these two locations was also measured, and the measurements combined with an approximate geometric relationships to calculate the Earth's circumference. Posidonius calculated the difference in the zenith angles at Canopus as about 1/48th of a circle between Rhodes and Alexandria. By estimating these two towns to be about 800 kilometers apart, he calculated the circumference of the Earth to be 38,600 kilometers. Again there were compensating errors, resulting in an accurate value. Another Greek scientist determined the circumference to be 28,960 kilometers, and unfortunately this shorter measurement was adopted by Ptolemy for his world maps. This estimate was widely accepted until the 1500s, when Gerardus Mercator revised the figure upward.

During the 17th and 18th centuries two developments led to intense activity directed

at measuring the size and shape of the Earth. Sir Isaac Newton and others reasoned the Earth must be flattened somewhat due to rotational forces. They argued that centrifugal forces cause the equatorial regions of the Earth to bulge as it spins on its axis. They proposed the Earth would be better modeled by an *ellipsoid*, a sphere that was slightly flattened at the poles. Measurements by their French contemporaries taken north and south of Paris suggested the Earth was flattened in an equatorial direction and not in a polar direction. The controversy persisted until expeditions by the French Royal Academy of Sciences between 1730 and 1745 measured the shape of the Earth near the equator in South America and in the high northern latitudes of Europe. Complex, repeated, and highly accurate measurements established that the curvature of the Earth was greater at the equator than the poles, and that an ellipsoid flattened at the poles was indeed the best geometric model of the Earth's surface.

Note that the words spheroid and ellipsoid are often used interchangeably. For example, the Clarke 1880 ellipsoid is often referred to as the Clarke 1880 spheroid, even though Clarke provided parameters for an ellipsoidal model of the Earth's shape. GIS

software often prompts the user for a spheroid when defining a coordinate projection, and then lists a set of ellipsoids for choices.

An ellipsoid is sometimes referred to as a special class of spheroid known as an "oblate" spheroid. Thus, it is less precise but still correct to refer to an ellipsoid more generally as a spheroid. It would perhaps cause less confusion if the terms were used more consistently, but the usage is widespread.

Specifying the Ellipsoid

Once the general shape of the Earth was determined, geodesists focused on precisely measuring the size of the ellipsoid. The ellipsoid has two characteristic dimensions (Figure 3-4). These are the *semi-major axis*, the radius a in the equatorial direction, and the *semi-minor axis*, the radius b in the polar direction. The equatorial radius is always greater than the polar radius for the Earth ellipsoid. This difference in polar and equatorial radii can also be described by the flattening factor, as shown in Figure 3-4.

Earth radii have been determined since the 18th century using a number of methods. The most common methods until recently have involved astronomical observations similar to the those performed by Posidonius. These astronomical observations, also called celestial observations, are combined with long-distance surveys over large areas (Figure 3-5). The distance and associated angles are measured in polar and equatorial directions, and used to estimate radii along the arcs. Several measurements were often combined to estimate semi-major and semi-minor axes.

Star and sun locations have been observed and cataloged for centuries, and combined with accurate clocks, the positions of these celestial bodies may be measured to precisely establish the latitudes and longitudes of points on the surface of the Earth. Measurements during the 18th, 19th and early 20th centuries used optical instruments for celestial observations (Figure 3-6).

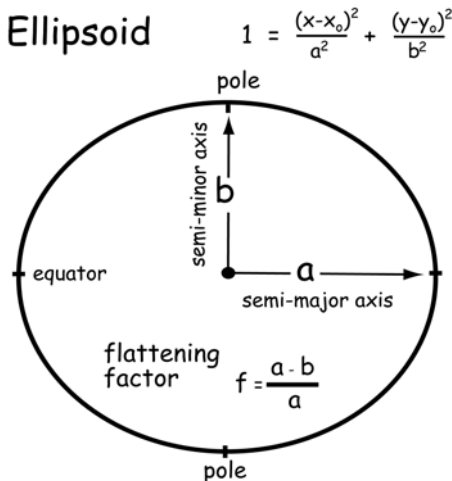
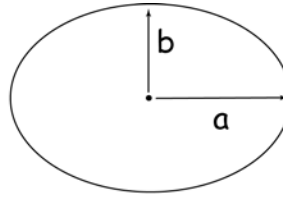


Figure 3-4: An ellipsoidal model of the Earth's shape.

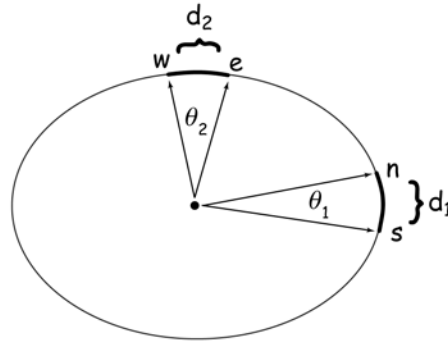
An ellipsoid is defined in part by two radii, a and b



We may use the relationship $d = r \cdot \theta$ to estimate radii:

$$a = \frac{d_1}{\theta_1}$$

$$b = \frac{d_2}{\theta_2}$$



Generally, the measurements are not at the poles and equator, and the math is more complicated, but the principle is the same.



Figure 3-5: Two arcs illustrate the surface measurements and calculations used to estimate the semi-major and semi-minor axes, here for North America. The arc lengths may be measured by surface surveys, and the angles from astronomical observations, as illustrated in Figure 3-2 and Figure 3-3.

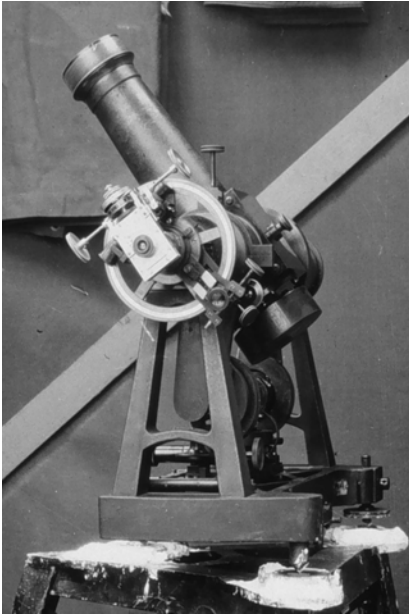


Figure 3-6: An instrument used in the early 1900s for measuring the position of celestial bodies.

Measurement efforts through the 19th and 20th centuries led to the establishment of a set of official ellipsoids (Table 3-1). Why not use the same ellipsoid everywhere on Earth, instead of the different ellipsoids listed in Table 3-1? The radii a and b for North America illustrated in Figure 3-5 would yield different estimates than those made in Europe or Africa, and there was no practical way to combine the measurements.

Historically, geodetic surveys were isolated by large water bodies. For example, surveys in Australia did not span the Pacific Ocean to reach Asia. Geodetic surveys relied primarily on optical instruments prior to the early 20th century. These instruments were essentially precise telescopes, and sighting distances were limited by the Earth's curvature. Individual survey legs greater than 50 kilometers (30 miles) were rare, so during this period there were no good ways to connect surveys between continents.

Table 3-1: Official ellipsoids. Radii may be specified more precisely than the 0.1 meter shown here (from Snyder, 1987 and other sources).

Name	Year	Equatorial Radius, a meters	Polar Radius, b meters	Flatten- ing Factor	Users
Airy	1830	6,377,563.4	6,356,256.9	1/299.32	Great Britain
Bessel	1841	6,377,397.2	6,356,079.0	1/299.15	Central Europe, Chile, Indonesia, U.S.
Clarke	1866	6,378,206.4	6,356,583.8	1/294.98	North America; Philip- pines
Clarke	1880	6,378,249.1	6,356,514.9	1/293.46	Most of Africa; France
Interna- tional	1924	6,378,388.0	6,356,911.9	1/297.00	Much of the world
Australian	1965	6,378,160.0	6,356,774.7	1/298.25	Australia
WGS72	1972	6,378,135.0	6,356,750.5	1/298.26	NASA, US Def. Dept.
GRS80	1980	6,378,137.0	6,356,752.3	1/298.26	Worldwide
WGS84	1987 - current	6,378,137.0	6,356,752.3	1/298.26	US DOD, Worldwide

Because continental surveys were isolated, ellipsoidal parameters were fit for each country, continent, or comparably large survey area. These ellipsoids represented continental measurements and conditions. Because of measurement errors, differences in methods for ellipsoidal calculation, and because the Earth's shape is not a perfect ellipsoid (described in the next section), different ellipsoids around the world usually had slightly different origins, axis orientations, and radii. These differences, while small, often result in quite different estimates for coordinate location at any given point, depending on the ellipsoid used.

More recently, data derived from satellites, lasers, and broadcast timing signals have been used for extremely precise measurements of relative positions across continents and oceans. Global measurements and faster computers allow us to estimate globally-applicable ellipsoids. These ellipsoids provide a “best” overall fit ellipsoid to observed measurements across the globe. Global ellipsoids such as the GRS80 or WGS84 are now preferred and most widely used.

The Geoid

As noted in the previous section, the true shape of the Earth varies slightly from the mathematically smooth surface of an ellipsoid. Differences in the density of the Earth cause variation in the strength of the gravitational pull, in turn causing regions to dip or bulge above or below a reference ellipsoid (Figure 3-7). This undulating shape is called a *geoid*.

Geodesists have defined the geoid as the three-dimensional surface along which the pull of gravity is a specified constant. The geoidal surface may be thought of as an imaginary sea that covers the entire Earth and is not affected by wind, waves, the Moon, or forces other than Earth's gravity. The surface of the geoid extends across the Earth, approximately at mean sea level across the oceans, and continuing under con-

tinents at a level set by gravity. The surface is always at right angles to the direction of local gravity, and this surface is the reference against which heights are measured.

Figure 3-8 shows how differences in the Earth's shape due to geoidal deviations will produce different best local ellipsoids. Surveys of one portion of the Earth that best fit the surveyed points will produce different best estimates of the ellipsoid origin, axis orientation, and of a and b than surveys of other parts of the Earth. Measurements based on Australian surveys yielded a different “best” ellipsoid than those in Europe. Likewise, Europe's best ellipsoidal estimate was different from Asia's, and from South America's, North America's, or those of other regions. One ellipsoid could not be fit to all the world's survey data because during the 18th and 19th centuries there was no clear way to combine a global set of measurements.

We must emphasize that a geoidal surface differs from mean sea level. Mean sea level may be higher or lower than a geoidal surface because ocean currents, temperature,

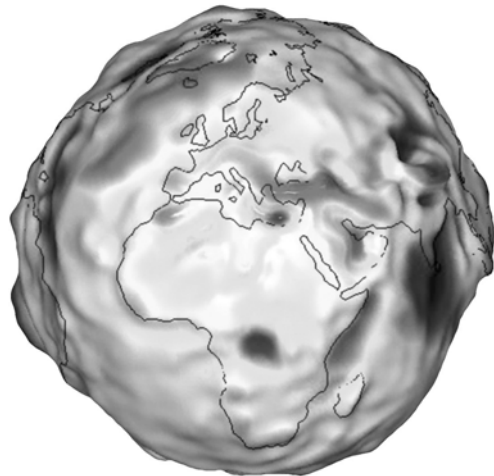


Figure 3-7: Depictions of the Earth's gravity field, as estimated from satellite measurements. These show the undulations, greatly exaggerated, in the Earth's gravity, and hence the geoid (courtesy University of Texas Center for Space Research, and NASA).

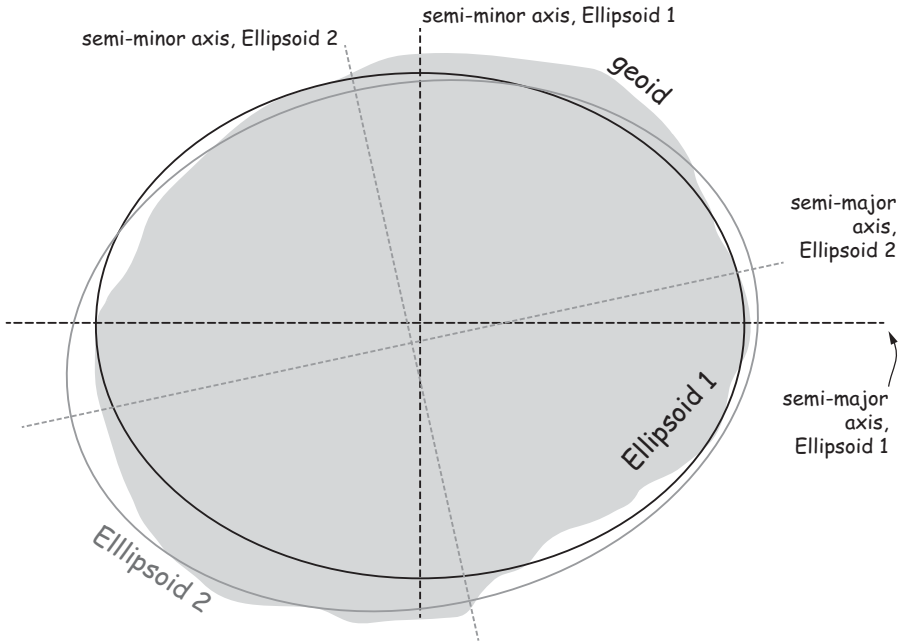


Figure 3-8: Different ellipsoids were estimated due to local irregularities in the Earth’s shape. Local best-fit ellipsoids varied from the global best fit, but until the 1970s, there were few good ways to combine global geodetic measurements.

salinity, and wind variations can cause persistent high or low areas in the ocean. These differences are measurable, in places over a meter (3 feet), perhaps small on global scale, but large in local or regional analysis. We historically referenced heights to mean sea level, and many believe we still do, but this is no longer true for most spatial data analyses.

Because we have two reference surfaces, a geoid and an ellipsoid, we also have two bases from which to measure height. Elevation is typically defined as the distance above a geoid. This height above a geoid is also called the *orthometric height* (Figure 3-9). Heights above an ellipsoid are often referred to as *ellipsoidal height*. These are illustrated in Figure 3-9, with the ellipsoidal height labeled *h*, and orthometric height labeled *H*. The difference between the ellipsoidal height and geoidal height at any location, shown in Figure 3-9 as *N*, has various names, including *geoidal height* and *geoidal separation*.

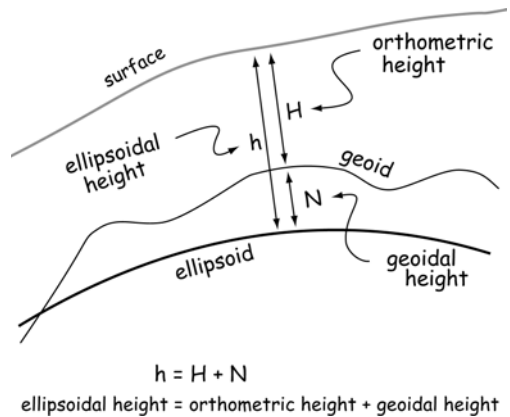


Figure 3-9: Ellipsoidal, orthometric, and geoidal height are interrelated. Note that values for *N* are highly exaggerated in this figure - values for *N* are typically much less than *H*.

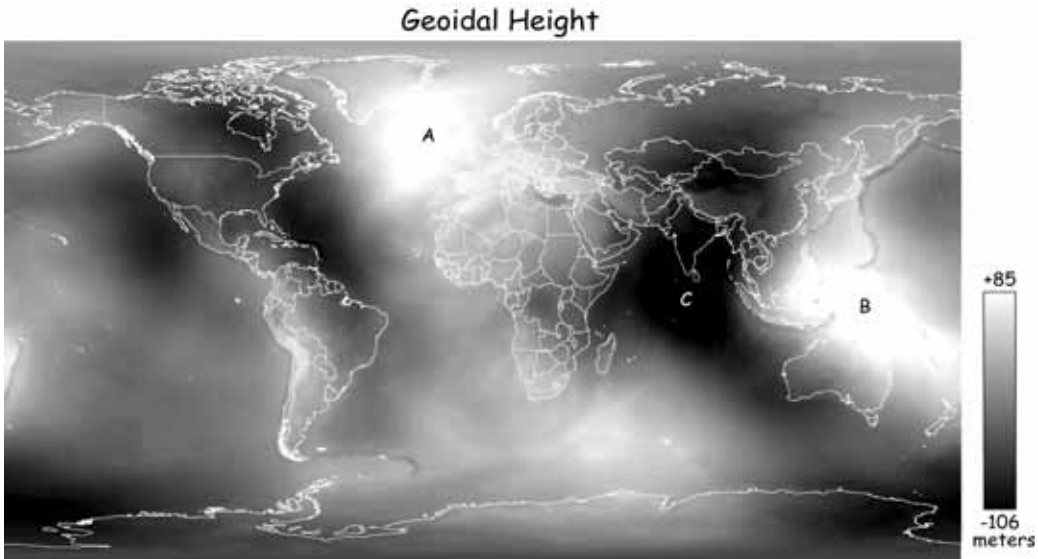


Figure 3-10: Geoidal heights vary across the globe. This figure depicts positive geoidal heights in lighter tones (geoid above the ellipsoid) and negative geoidal heights in darker tones. Note that geoidal heights are positive for large areas near Iceland and the Philippines (A and B, respectively), while large negative values are found south of India (C). Continental and country borders are shown in white.

The absolute value of the geoidal height is less than 100 meters over most of the Earth (Figure 3-10). Although it may at first seem difficult to believe, the “average” ocean surface near Iceland is more than 150 meters “higher” than the ocean surface northeast of Jamaica. This height difference is measured relative to the ellipsoid. Since gravity pulls in a direction that is perpendicular to the geoidal surface, the force is at a right angle to the surface of the ocean, resulting in permanent bulges and dips in the mean ocean surface due to variations in the gravitational pull. Variation in ocean heights due to swells and wind-driven waves are more apparent at local scales, but are much smaller than the long-distance geoidal undulations.

The geoidal height is quite small relative to the polar and equatorial radii. As noted in Table 3-1, the Earth’s equatorial radius is about 6,780,000 meters, or about 32,000 times the range of the highest to lowest geoidal heights. This small geoidal height is imperceptible in an object at human scales. For example, the largest geoidal height is less than the relative thickness of a coat of

paint on a ball three meters (10 feet) in diameter. However, while relatively small, the geoidal variations in shape must still be considered for accurate vertical and horizontal mapping over continental or global distances.

The geoid is a measured and interpolated surface, and not a mathematically defined surface. The geoid’s surface is measured using a number of methods, initially by a combination of *plumb bob*, a weight suspended by a string that indicates the direction of gravity, and horizontal and vertical distance measurements, and later with various types of *gravimeters*, devices that measure the gravitational force.

Satellite-based measurements in the late 20th century substantially improved the global coverage, quality, and density of geoidal height measurements. The GRACE experiment, initiated with the launch of twin satellites in 2002, is an example of such improvements. Distances between a pair of satellites are constantly measured as they orbit the Earth. The satellites are pulled closer or drift farther from the Earth due to variation in the gravity field. Because the

orbital path changes slightly each day, we eventually have nearly complete Earth coverage of the strength of gravity, and hence the location of the reference gravitational surface. The ESA GOCE satellite, launched in 2009, uses precision accelerometers to measure gravity-induced velocity change. GRACE and GOCE observations have substantially improved our estimates of the gravitational field and geoidal shape.

Satellite and other observations are used by geodesists to develop geoidal models. These support a series of geoid estimates, e.g., by the U.S. NGS with GEOID90 in 1990, with succeeding geoid estimates in 1993, 1996, 1999, 2003, 2009, and one planned for 2012. These are called models because we measured geoidal heights at points or along lines at various parts of the globe, but we need geoidal heights everywhere. Equations are statistically fit that relate the measured geoidal heights to geographic coordinates. Given any set of geographic coordinates, we may then predict the geoidal height. These models provide an accurate estimation of the geoidal heights for the entire globe.

Geographic Coordinates, Latitude, and Longitude

Once a size and shape of the reference ellipsoid has been determined, the Earth poles and equator are also defined. The poles are defined by the axis of revolution of the ellipsoid, and the equator is defined as the circle mid-way between the two poles, at a right angle to the polar axis, and spanning the widest dimension of the ellipsoid. We estimate these locations from precise surface and astronomical measurements. Once the locations of the polar axis and equator have been estimated, we can define a set of geographic coordinates. This creates a reference system by which we may specify the position of features on the ellipsoidal surface.

As noted in Chapter 2, geographic coordinate systems consist of latitude, which varies from north to south, and longitude, which varies from east to west (Fig-

ure 3-11). Lines of constant longitude are called meridians, and lines of constant latitude are called parallels. Parallels run parallel to each other in an east-west direction around the Earth. The meridians are geographic north/south lines that converge at the poles.

By convention, the equator is taken as zero degrees latitude, and latitudes increase to maximum values of 90 degrees in the north and south. Latitudes are thus designated by their magnitude and direction, for example 35°N or 72°S. When signed values are required, northern latitudes are designated positive and southern latitudes are designated negative. An international meeting in 1884 established a longitudinal origin intersecting the Royal Greenwich Observatory in England. Known as the *prime* or *Greenwich meridian*, this north-to-south line was the origin, or zero value, for longitudes. Improvements in measurements, crustal movements, and changes in conventions have resulted in the present zero longitude about 102 meters (335 feet) east of the Greenwich observatory. East or west longitudes are specified as angles of rotation away from the Prime Meridian, from -180 (westerly) to +180 (easterly).

There is often confusion between magnetic north and geographic north. Magnetic

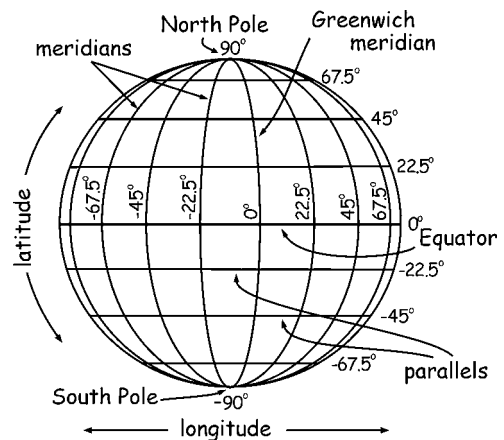


Figure 3-11: Nomenclature of geographic latitudes and longitudes.

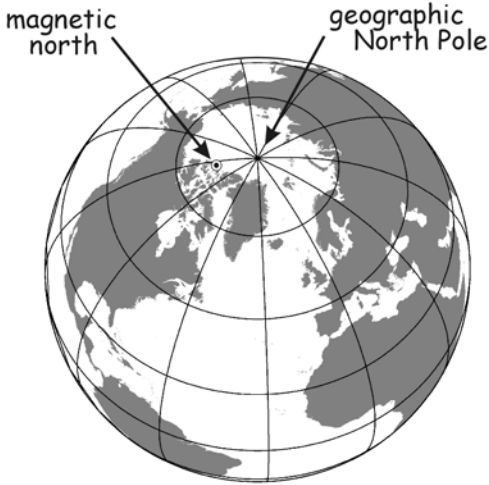


Figure 3-12: Magnetic north and the geographic North Pole.

north and the geographic north do not coincide (Figure 3-12). Magnetic north is the location towards which a compass points. The geographic North Pole is the northern pole of the Earth’s axis of rotation. If you were standing on the geographic North Pole with a compass, it would point approximately in the direction of northern Canada, towards magnetic north some 600 kilometers away.

Because magnetic north and the geographic North Pole are not in the same place, a compass does not point at geographic north when observed from most places on Earth.

The compass will usually point east or west of geographic north, defining an angular difference in direction to the poles. This angular difference is called the magnetic *declination* and varies across the globe. The specification of map projections and coordinate systems is always in reference to the geographic North Pole, not magnetic north.

Geographic coordinates do not form a Cartesian system (Figure 3-13). A Cartesian system defines lines of equal value in a right-angle grid. Geographic coordinates occur on a curved surface, and the longitudinal lines converge at the poles. This convergence means the distance spanned by a degree of longitude varies from south to north. A degree of longitude spans approximately 111.3 kilometers at the equator, but 0 kilometers at the poles. In contrast, the ground distance for a degree of latitude varies only slightly, from 110.6 kilometers at the equator to 111.7 kilometers at the poles.

Convergence causes regular geometric figures specified in geographic coordinates to appear distorted when drawn on a globe (Figure 3-13, left). For example, “circles” with a fixed radius in geographic units, such as 5°, are not circles on the surface of the globe, although they may appear as circles when the Earth surface is “unrolled” and plotted with distortion on a flat map; note the erroneous size and shape of Antarctica at the bottom of Figure 3-13.

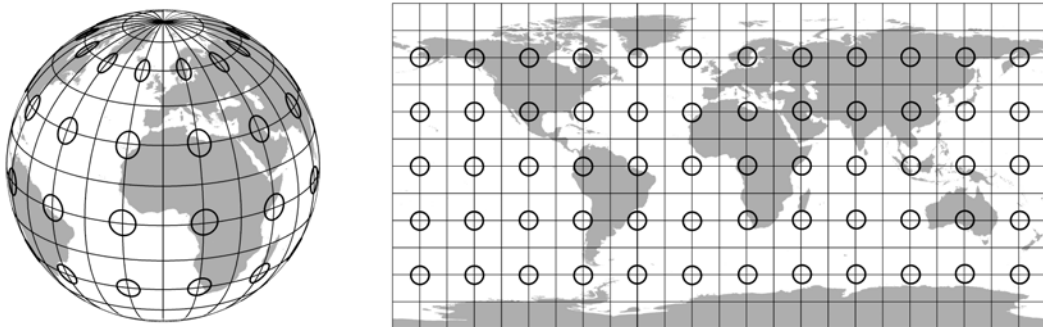


Figure 3-13: Geographic coordinates on a spherical (left) and Cartesian (right) representation. Notice the circles with a 5 degree radius appear distorted on the spherical representation, illustrating the change in surface distance represented by a degree of longitude from the equator to near the poles.

Horizontal Datums

The geographic coordinate system we have just described provides for specifying locations on the Earth. However, this gives us the exact longitude of only one arc, the zero line of longitude. We must estimate the longitudes and latitudes of all other locations through surveying measurements; until recently by observing stars and by measuring distances and directions between points. These surveying methods have since been replaced by modern, satellite-based positioning, but even these new methods are ultimately dependent on astronomical observations. Through these methods we establish a set of points on Earth for which the horizontal and vertical positions have been accurately determined.

These well-surveyed points allow us to specify a *reference frame*, including an origin or starting point. If we are using a spherical reference frame, we must also specify the orientation and scale of our ellipsoid. If we are using a three-dimensional Cartesian reference frame, we must specify the X, Y, and Z axes, including their origin and orientation. All other coordinate locations we use are measured with reference to this set of precisely surveyed points, including the coordinates we enter in our GIS to represent spatial features.

Many countries have a government body charged with making precise geodetic surveys. For example, most surveys in the United States are related back to high accuracy points maintained by the National Geodetic Survey (NGS). The NGS establishes geodetic latitudes and longitudes of known points, most of which are monumented with a bronze disk, concrete posts, or other durable markers. These points, taken together, underpin *geodetic datums*, upon which most subsequent surveys and positional measurements are based.

A *datum* is a reference surface. A geodetic datum consists of two major components. The first component is an ellipsoid with a spherical or three-dimensional Cartesian coordinate system and an origin. Eight

parameters are needed to specify the ellipsoid: a and b to define the size/shape of the ellipsoid, the X, Y, and Z values of the origin, and an orientation angle for each of the three axes.

The second part of a useful datum consists of a set of points and lines that have been painstakingly surveyed using the best methods and equipment, and an estimate of the coordinate location of each point in the datum, e.g., the NGS points described in the previous paragraphs. Some authors define the datum as a specified reference surface, and a *realization of a datum* as that surface plus a physical network of precisely measured points. In this nomenclature, the measured points describe a *Terrestrial Reference Frame*. This clearly separates the theoretical surface, the reference system or datum, from the terrestrial reference frame, a specific set of measurement points that help fix the datum. While this more precise language may avoid some confusion, datum will con-

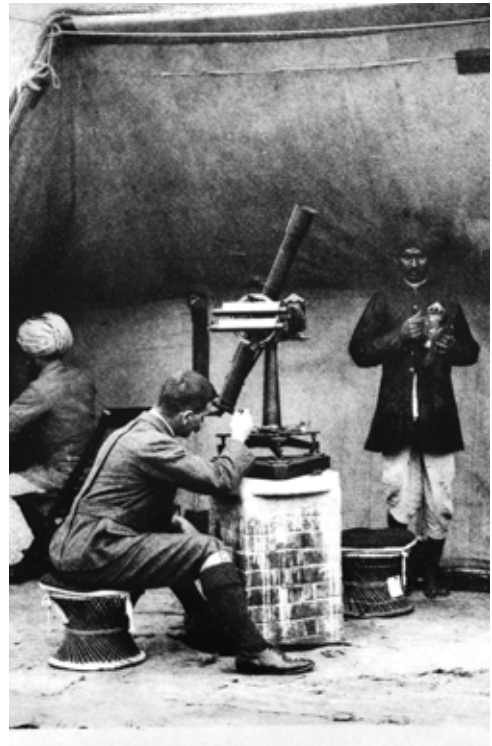


Figure 3-14: Astronomical observations were used in early geodetic surveys to measure datum locations (courtesy NMSI)

tinue to refer to both the defined surface and the various realizations of each datum.

Different datums are specified through time because our realizations, or estimates of the datum, change through time. New points are added and survey methods improve. We periodically update our datum when a sufficiently large number of new survey points has been measured. We do this by re-estimating the coordinates of our datum points after including these newer measurements, thereby improving our estimate of the position of each point.

Most governments have an official body responsible for measuring and maintaining a set of datum points. They keep records of the points locations, estimate new datums, and distribute datum descriptions and point coordinate values. Without these, surveyors, GIS practitioners, and others cannot precisely identify coordinate location.

Precisely surveyed points are commonly known as *bench marks*. Bench marks usually consist of a brass disk embedded in rock or concrete (Figure 3-15), although they also may consist of marks chiseled in rocks, embedded iron posts, or other long-term marks. Due to the considerable effort and cost of establishing the coordinates for each bench mark, they are often redundantly



Figure 3-15: A brass disk used to monument a survey bench mark.



Figure 3-16: Signs are often placed near control points to warn of their presence and aid in their location.

monumented, and their distance and direction from specific local features are recorded. Control survey points are often identified with a number of nearby signs to aid in recovery (Figure 3-16).

Geodetic surveys in the 18th and 19th centuries combined horizontal measurements with repeated, excruciatingly precise astronomical observations to determine latitude and longitude of a small set of points. Only a few datum points were determined using astronomical observations. Astronomical observations were typically used at the starting point, a few intermediate points, and near the end of geodetic surveys. This is because star positions required repeated measurements over several nights. Clouds, haze, or a full moon often lengthened the measurement times. In addition, celestial measurements required correction for atmospheric refraction, a process which bends light and changes the apparent position of stars. Refraction depends on how high the star is in the sky at the time of measurement, as well as temperature, atmospheric humidity, and other factors.

Historically, horizontal measurements were as precise and much faster than astronomical measurements when surveying over counties or state-sized regions. These horizontal surface measurements were used to

connect astronomically surveyed points and thereby create an expanded, well-distributed set of known datum points. Figure 3-17 shows an example survey, where open circles signify points established by astronomical measurements and filled circles denote points established by surface measurements.

Figure 3-17 shows a *triangulation survey*, until the mid 1980s (and the advent of GPS) the method commonly used to establish datum points via horizontal surface measurements. Triangulation surveys utilize a network of interlocking triangles to determine positions at survey stations. Triangulation surveys were adopted because we can create them through angle measurement, with few surface distance measurements, an advantage in the late 18th century when many datums were first developed. Triangulation also improves accuracy; because there are multiple measurements to each survey station, the location at each station may be computed by various paths. The survey accuracy can be field-checked, because large

differences in a calculated station location via different paths indicate a survey error. There are always some differences in the measured locations when traversing different paths. An acceptable error limit was often set, usually as a proportion of the distance surveyed. In one common standard, differences in the measured location of more than 1 part in 100,000 would be considered unacceptable. When unacceptable errors were found, survey lines were re-measured.

Triangulation networks spanned long distances, from countries to continents (Figure 3-18). Individual measurements of these triangulation surveys were rarely longer than a few kilometers, however triangulations were nested, in that triangulation legs were combined to form larger triangles spanning hundreds of kilometers. These are demonstrated in Figure 3-18 where the sides of each large triangle are made up themselves of smaller triangulation traverses.

Datum Adjustment

Once a sufficiently large set of points have been surveyed, the survey measurements must be harmonized into a consistent set of coordinates. Small inconsistencies are inevitable in any large set of measurements, causing ambiguity in locations. In addition, historically the long reaches spanned by the triangulation networks, as shown in Figure 3-18, could be helpful in recalculating certain constants, such as the Earth's curvature (Figure 3-5). Later, satellite-based measurements were used to better estimate other constants, such as the datum origin. The positions of all points in a reference datum are estimated in a network-wide *datum adjustment*. The datum adjustment reconciles errors across the network, first by weeding out blunders or obvious mis-measurements or other mistakes, and also by mathematically minimizing errors by combining repeat measurements and statistically assigning higher influence to consistent or more precise measurements. Note that a given datum adjustment only incorporates measurements up to a given point in time,

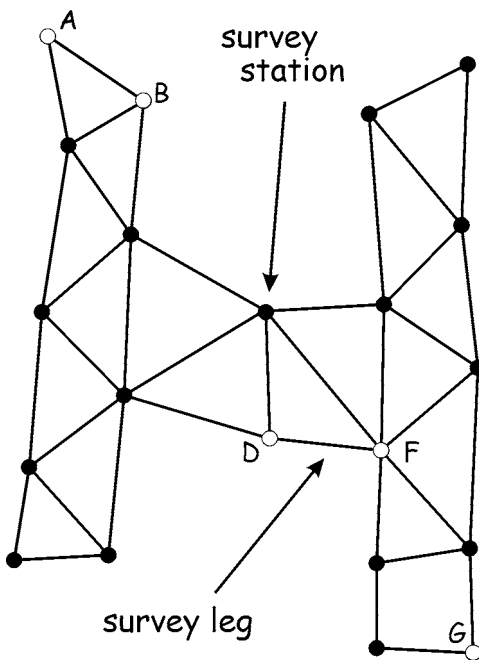


Figure 3-17: A triangulation survey network. Stations may be measured using astronomical (open circles) or surface surveys (filled circles).

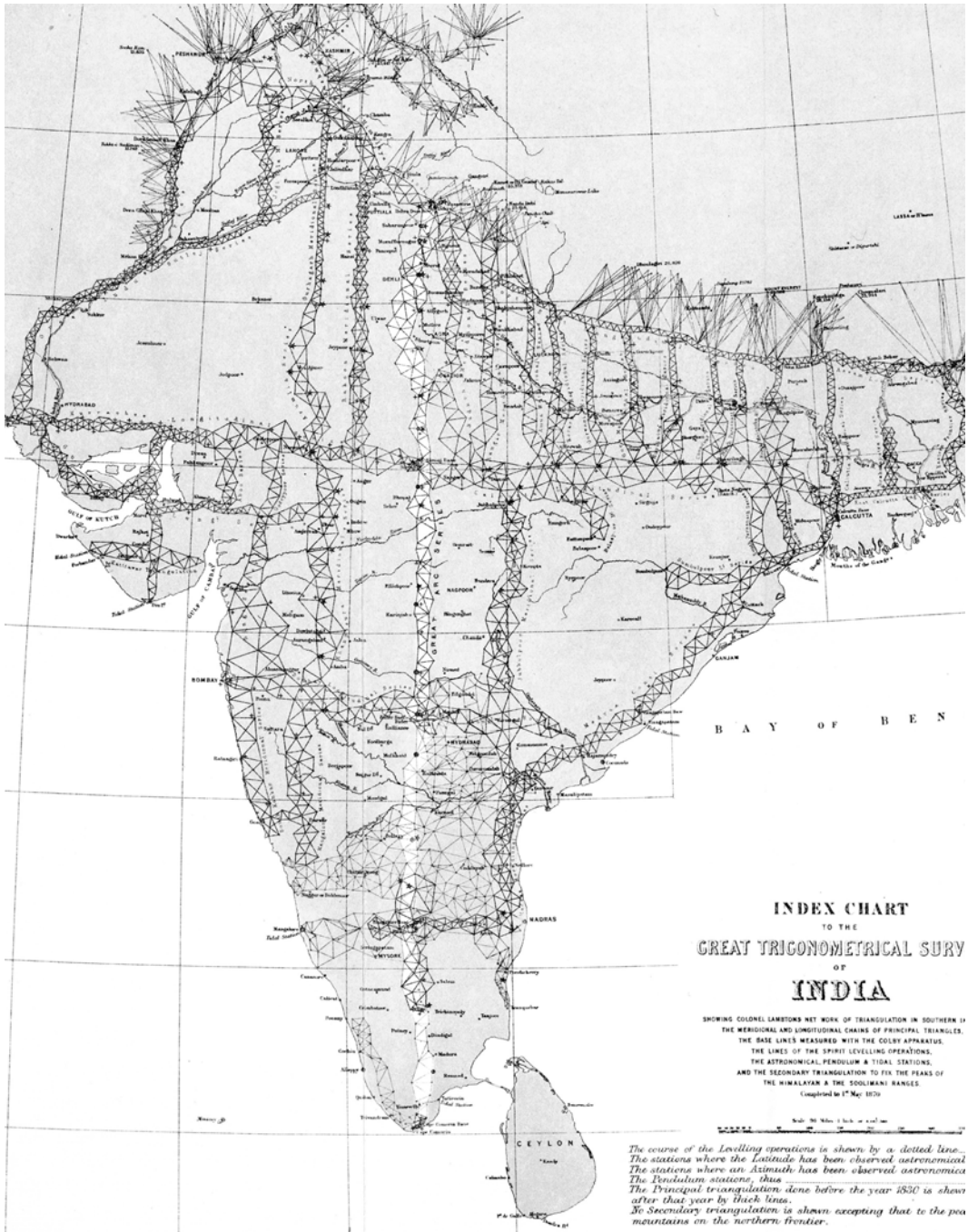


Figure 3-18: A map of the triangulation survey network established across India in the 1800s. Each leg of the triangles, shown here as a single line, is in turn a triangulation survey. This nested triangulation provides reinforcing measurements, thereby increasing the accuracy of the surveyed positions (courtesy NMSI).

and may be viewed as our best estimate, at that point, of the measured set of locations.

Periodic datum adjustments result in series of regional or global reference datums. Each datum is succeeded by an improved, more accurate datum. The calculation of a new datum requires that all surveys must be simultaneously adjusted to reflect our current “best” estimate of the true positions of each datum point. Generally a statistical least-squares adjustment is performed, but this is not a trivial exercise, considering the adjustment may include survey data for tens of thousands of old and newly surveyed points from across the continent, or even the globe. Because of their complexity, these continent-wide or global datum calculations have historically been quite infrequent. Computational barriers to datum adjustments have diminished in the past few decades, and so datum adjustments and new versions of datums are now more frequent.

A datum adjustment usually results in a change in the coordinates for all existing datum points, as coordinate locations are estimated for both old and new datum points. The datum points do not move, but

our best estimates of the datum point coordinates will change. Differences between the datums reflect differences in the control points, survey methods, and mathematical models and assumptions used in the datum adjustment.

Figure 3-19 illustrates how ellipsoids might change over time, even for the same survey region. Ellipsoid A is estimated with the datum coordinates for pt1 and pt2, with the shown corresponding coordinate axes, origin, and orientation. Ellipsoid B is subsequently fit, after pts 3 through 7 have been collected. This newer ellipsoid has a different origin and orientation for its axis, causing the coordinates for pt1 and pt2 to change. The points have not moved, but the best estimate of their locations, relative to the origin set by the new, more complete set of datum points, will have changed. You can visualize how the latitude angle from the origin to pt1 will change because the origin for ellipsoid A is in a different location than the origin for ellipsoid B. This apparent, but not real, movement is called the datum shift, and is expected with datum adjustments.

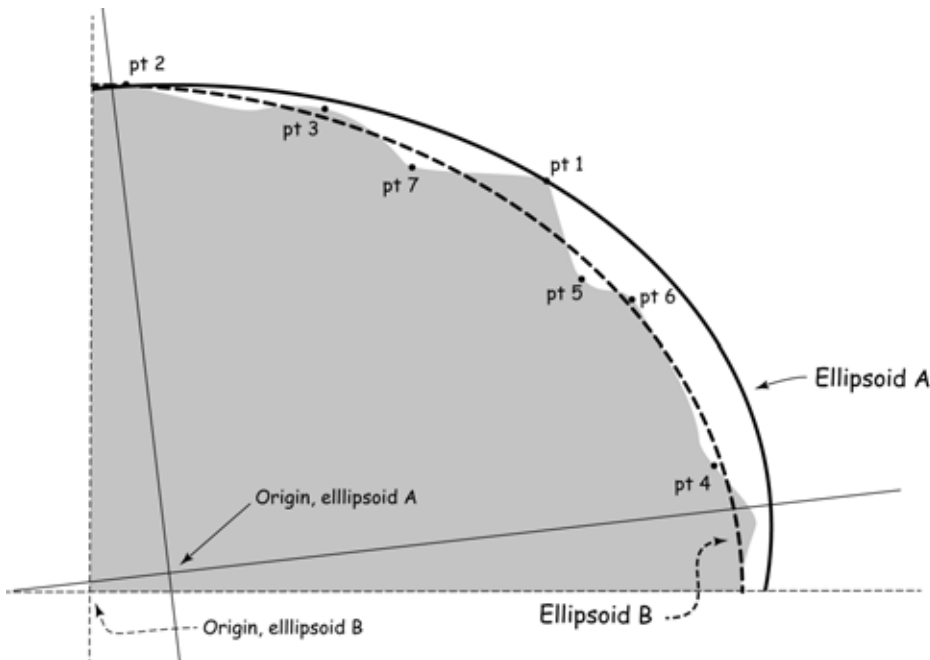


Figure 3-19: An illustration of two datums, one corresponding to Ellipsoid A and based on the fit to pt1 and pt2, and a subsequent datum resulting in Ellipsoid B, and based on a fit of pt1 through pt7.

Commonly Used Datums

Three main series of horizontal datums have been used widely in North America. The first of these is the *North American Datum of 1927* (NAD27). NAD27 is a general least-squares adjustment that included all horizontal geodetic surveys completed at that time. The geodesists used the Clarke Ellipsoid of 1866 and held fixed the latitude and longitude of a survey station in Kansas. NAD27 yielded adjusted latitudes and longitudes for approximately 26,000 survey stations in the United States and Canada.

The *North American Datum of 1983* (NAD83) is the successor datum to NAD27. We place the (1986) after the NAD83 designator to indicate the year, or version, of the datum adjustment. It was undertaken by the NGS to include the large number of geodetic survey points established between the mid-1920s and the early 1980s. Approximately 250,000 stations and 2,000,000 distance measurements were included in the adjustment. The GRS80 ellipsoid was used, and NAD83(1986) is Earth-centered reference, rather than fixing a station as with NAD27. The shifts in estimated coordinate locations between NAD27 and NAD83(1986) were large, on the order of 10's to up to 200 meters in North America. In most instances the surveyed points physically moved very little, e.g., due to tectonic shifts, but our best estimates of point location changed by as much as 200 meters.

Precise GPS data became widely available soon after the initial NAD83(1986) adjustment, and these were often more accurate than NAD83(1986) position estimates. Between 1989 and 2004, the NGS collaborated with other federal agencies, state and local governments, and private surveyors in creating *High Accuracy Reference Networks* (HARNs), also known as *High Precision Geodetic Networks* (HPGN) in each state and most U.S. territories.

Subsequent NAD83 adjustments have incorporated measurements from the Continuously Operating Reference Station (CORS) network (Figure 3-20). This growing net-

work of satellite observation stations allowed improved datum realizations, including NAD83(CORS93), NAD83(CORS94), NAD83(CORS96), NAD83(2007), and NAD83(2011). The NAD83(2007) datum may be viewed as a successor to the NAD83(HARN). Approximately 70,000 high-accuracy GPS points were adjusted with reference to the NAD83(CORS96) coordinates for the CORS network. NAD83(2011) is a long-observation adjustment based on CORS stations, with coordinates re-estimated for a broad set of bench marks. This datum realization allows surveyors to obtain the coordinates for a widespread set of physical locations, which may then be used as a starting point for subsequent surveys.

Position estimates of locations change by a few centimeters when compared among the NAD83(CORSxx) datums, important improvements for geodesists and extremely precise surveying, but small relative to spatial error budgets for many GIS projects. Differences among current and future NAD83(CORSxx) datums are likely to remain small, on the order of a few centimeters or less in tectonically stable areas, as newer NAD83 datum adjustments are calculated in the future.

The *World Geodetic System of 1984* (WGS84) is a set of datums developed and primarily used by the U.S. Department of Defense (DOD). It was introduced in 1987 based on Doppler satellite measurements of the Earth, and is used in most DOD maps and positional data. The WGS84 ellipsoid is similar to the GRS80 ellipsoid. WGS84 has been updated with more recent satellite measurements and is specified using a version designator. The update based on data collected up to January 1994 is designated as WGS84 (G730). WGS84 datums are not widely used outside of the military because they are not tied to a set of broadly accessible, documented physical points.

There have been several subsequent WGS84 datum realizations. The original datum realization exhibited positional accuracy of key datum parameters to within

between one and two meters. Subsequent satellite observations improved accuracies. A re-analysis was conducted on data collected through week 730 of the GPS satellite schedule, resulting in the more accurate WGS84(G730). Successive re-adjustments in weeks 873 and 1150 are known as WGS84(G873) and WGS(G1150), respectively. There will likely be more adjustments in the future.

It has been widely stated that the original WGS84 and NAD83(86) datums were essentially equivalent. Both used the GRS80 ellipsoid, but the defining document for WGS84 notes differences of up to two meters between point locations measured against NAD83(86) versus the original WGS84 datum realizations. These differences have remained through time. You should note that there are positional differences among and between all versions of

both NAD83 and of WGS84, and ignoring the differences may result in positional error. The NAD83 and WGS84 datums have always been different by up to two meters.

Another set of datums used worldwide, known as the *International Terrestrial Reference Frames*, (ITRF), are realizations of the International Terrestrial Reference System (ITRS). A primary purpose for ITRS is to estimate continental drift and crustal deformation by measuring the location and velocity of points, using a worldwide network of measurement locations. Each realization is noted by the year, e.g., ITRF89, ITRF90, ITRF91, and so forth, and includes the X, Y, and Z location of each point and the velocity of each point in three dimensions. The European Terrestrial Reference System (ETRS89 and frequent updates thereafter) is based on ITRF measurements.

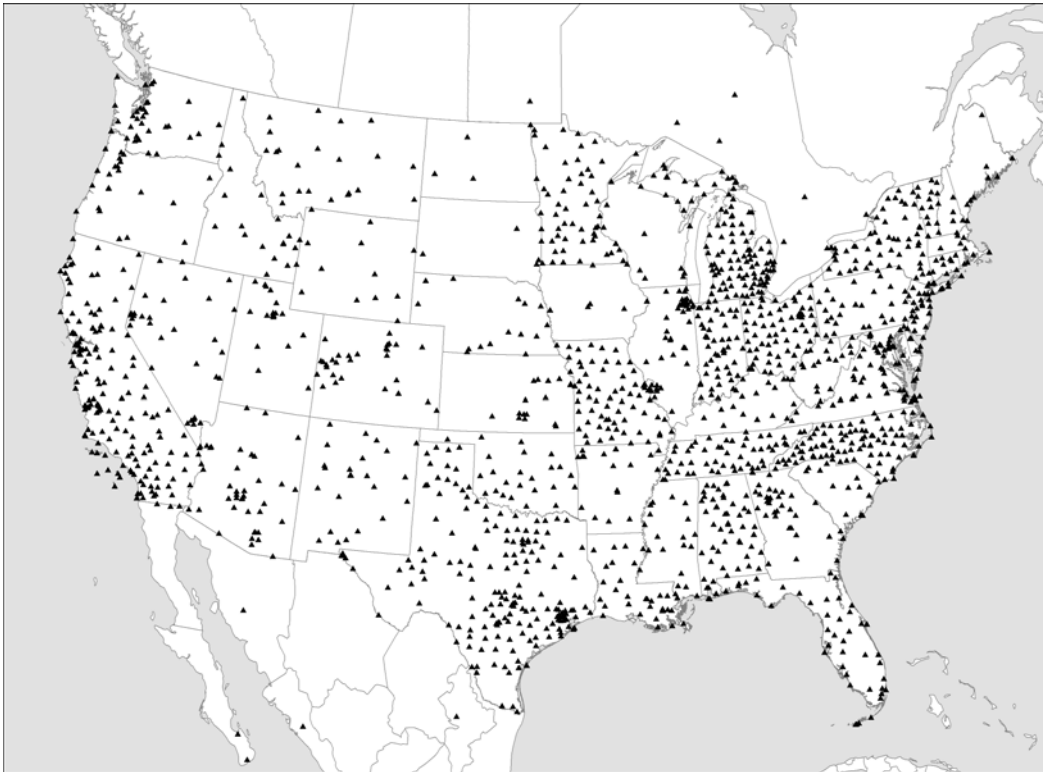


Figure 3-20: Partial distribution of the CORS network, as of 2008. This evolving network is the basis for the NAD83 (1993), NAD83(1994), NAD83(1996), and subsequent U.S. datum adjustments.

As noted earlier, different datums are based on different sets of measurements and ellipsoids, causing the coordinates for bench mark datum points to differ between datums and realizations. Differences are typically largest between legacy pre-satellite datum realizations, and post-satellite measurement datums. For example, the latitude and longitude location of a given bench mark in the NAD27 datum will likely be different from the latitude and longitude of that same bench mark in NAD83 or WGS84 datums by tens of meters, and up to 80 meters. This is described as a *datum shift*.

Figure 3-21 indicates the relative size of datum shifts at an NGS bench mark between NAD27 and NAD83(86) at one point in the eastern U.S., based on estimates provided by the National Geodetic Survey. Notice that the datum shift between NAD27 and

NAD83(86) is quite large, approximately 40 meters (140 feet), typical of the up to 100's of meters of shifts from pre-satellite, regional datums to post-satellite, global datums.

A datum shift does not imply that points have moved. Most monumented points are stationary relative to their immediate surroundings. The locations change over time as the large continental plates move, but these changes are small, on the order of a few millimeters per year, except in tectonically active areas such as coastal California; for most locations it is just our estimates of the coordinates that have changed. As survey measurements improve through time and there are more of them, we obtain better estimates of the true locations of the monumented datum points.

Examples of Datum Shifts

Successive datum transformations for New Jersey control point, Bloom 1

Datum	Longitude (W)	Latitude(N)	Shift(m)
NAD27	74° 12' 3.86927"	40° 47' 0.76531"	} 36.3
NAD83(1986)	74° 12' 2.39240"	40° 47' 1.12726"	
NAD83(HARN)	74° 12' 2.39069"	40° 47' 1.12762"	} 0.04
NAD83(CORS96)	74° 12' 2.39009"	40° 47' 1.12936"	} 0.05
NAD83(2007)	74° 12' 2.38977"	40° 47' 1.12912"	} 0.01
WGS84(G1150)	74° 12' 2.39720"	40° 47' 1.15946"	} 0.95

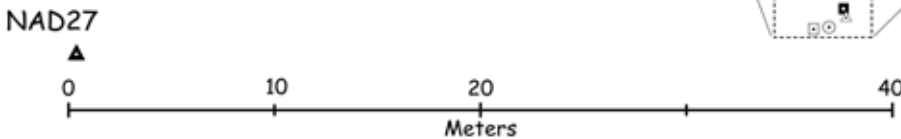
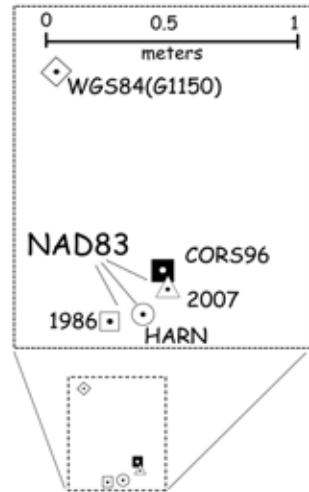


Figure 3-21: Datum shifts in the coordinates of a point for some common datums. Note that the estimate of coordinate position shifts approximately 36 meters from the NAD27 to the NAD83(1986) datum, while the shift from NAD83(1983) to NAD83(HARN) then to NAD83(CORS96) are approximately 0.05 meters. The shift to WGS84(G1150) is also shown, here approximately 0.95 m. Note that the point may not be moving, only our datum estimate of the point's coordinates. Calculations are based on NGS NAD-CON and HTDP software.

We must emphasize while much data are collected in WGS84 datums using GNSS (such as GPS), most data are converted to a local or national datum before use in a GIS. In the United States, this typically involves GNSS accuracy augmentation, often through a process called differential correction, described in detail in Chapter 5. Corrections are often based on an NAD83 datum, effectively converting the coordinates to the NAD83 reference, but ITRF datums are also commonly used. Ignorance of this “implicit” transformation is a common source of error in spatial data, and should be recognized.

For a datum to be practically useful in a GIS, we typically need the datum coordinates for a widely distributed and uniformly documented set of monumented bench marks. The development of new data through local surveys and image interpretation requires that we tie our new data to this existing network of surveyed points. In the U.S., most spatial data are tied to the widely distributed set of bench marked points reported in the NAD83 (CORSxx) datums, and state, county, and local surveys referenced to these points. The error introduced in ignoring the differences between versions of WGS84 and the NAD83 or other local datums can be quite large, generally up to 2 meters or more (Figure 3-21). Errors in ignoring differences among older datums are larger still, up to 100’s of meters. We must use a technique called a datum transformation to combine spatial data measured relative to different datums.

This conversion often happens implicitly when processing the GNSS data. As described in Chapter 6, most precise GNSS data results from correcting field measurements at unknown points against simultaneous GNSS measurements at a known point. This differential correction is usually configured such that the resulting coordinates are expressed in the same datum as the known points. If the correction sources are expressed in NAD83(CORS96) coordinates, corrected positions are initially created in these NAD83(CORS96) coordinates.

There are a few points about datums that must be emphasized. First, different datums mean different coordinate systems. You do not expect coordinates for any physical point to be the same when they are expressed relative to different datums.

Second, the version of the datum is important. NAD83(1986) is a different realization than NAD83(1996). The datum is incompletely specified unless the version is noted. Many GIS software packages refer to a datum without the version, e.g., NAD83. This is indeterminate, and confusing, and shouldn’t be practiced. It forces the user to work with ambiguity.

Third, differences between families of datums change through time. The NAD83(86) datum realization is up to two meters different than the NAD83(CORS96), and the original WGS84 differs from the current version by more than a meter over much of the Earth. Differences in datum realizations depend on the versions and location on Earth. This means you should assume all data should be converted to the same datum, via a datum transformation, before combination in a GIS. This rule may be relaxed if the datum difference errors are small compared to other sources of error, or to the data accuracy required for the intended spatial analysis.

Datum Transformations

Estimating the shift and converting geographic coordinates from one datum to another typically requires a *datum transformation*. A datum transformation provides the latitude and longitude of a point in one datum when we know them in another datum, for example, we can calculate the latitude and longitude of a bench mark in NAD83(HARN) when we know these geographic coordinates in NAD83(CORS96) (Figure 3-22).

Datum transformations are often more complicated when they involve older datums. Many older datums were created piecemeal to optimize fit for a country or continent. The amount of shift between one

datum and another often varies across the globe because the errors in measurements may be distributed idiosyncratically. Measurements in one area or period may have been particularly accurate, while in another area or time they may exhibit particularly large errors. Combining them in the datum adjustment affect the local and global differences among datums in their own unique way. Simple formulas often do not exist for transformations involving many older datums, for example from NAD27 to NAD83. Specialized datum transformations may be provided, usually by government agencies, using a number of different methods. As an example, in the United States the National Geodetic Survey has published a number of papers on datum transformations and provided datum transformation software tools, including NADCON to convert between NAD27 and NAD83 datums.

Transformation among newer datums may use more general analytical approaches that apply mathematical transformations between three-dimensional, Cartesian coordinate systems (Figure 3-22). These Earth or near-Earth centered (geocentric) coordinate systems allow conversion among most GPS and CORS-based NAD83, WGS84 and ITRF systems, and are supported in large part by improved global measurements from artificial satellites, as described in the previous few pages. This three-dimensional approach typically allows for a shift in the origin, a rotation, and a change in scale from one datum to another.

A mathematical geocentric datum transformation is typically a multi-step process. These datum transformations are based on one of a few methods, for example, in past times a *Molodenski transformation* using a system of equations with three or five parameters, or more currently, a *Helmert*

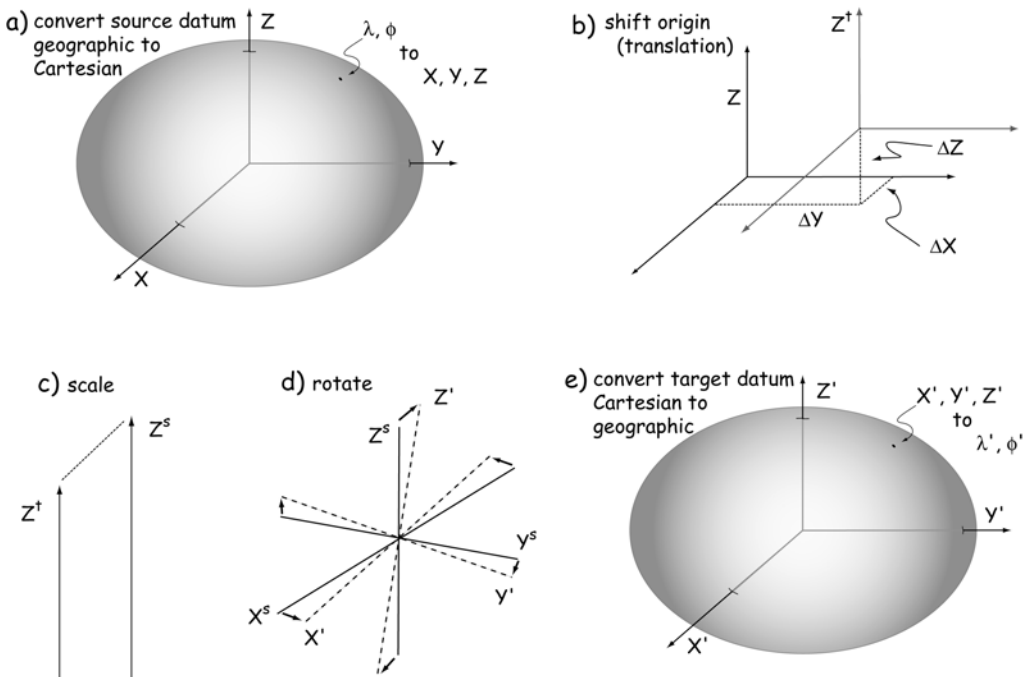


Figure 3-22: Application of a modern datum transformation. Geographic coordinates (longitude, λ , and latitude, ϕ), are transformed to a new datum by a) conversion from geographic to Cartesian coordinates in the old datum (through a set of equations that are not shown), b) applying an origin shift, c) scaling and d) rotating these shifted coordinates, and e) converting these target datum Cartesian coordinates, X', Y', Z' , to the longitude and latitude, λ', ϕ' , in the target datum.

transformation using seven parameters (Figure 3-22). First, geographic coordinates on the source datum are converged from longitude (λ) and latitude (ϕ) to X, Y, and Z Cartesian coordinates. An origin shift (translation), rotation, and scale are applied. This system produces new X', Y', and Z' coordinates in the target datum. These X', Y', and Z' Cartesian coordinates are then converted back to geographic coordinates, longitudes and latitudes (λ' and ϕ'), in the target datum.

More advanced methods allow these 7 transformation parameters to change through time, as tectonic plates shift, for a total of 14 parameters. These methods are incorporated into software that calculate transformations among modern datums, for example, the Horizontal Time Dependent Positioning (HTDP) tool available from the U.S. NGS (www.ngs.noaa.gov/TOOLS/Htdp/Htdp.shtml). HTDP converts among recent NAD83 datums and most ITRF and WGS84 datums.

Positions also change through time as tectonic plates shift, so that the most precise geodetic measurements refers to the epoch, or fixed time period, at which the point was measured. The HTDP software includes options to calculate the shift in a location due to measuring against different reference datums (e.g., NAD83(CORS96) to WGS84(G1150), the shift due to different realizations of a datum (e.g., NAD83(CORS96) to NAD83(2011)), the shift due to measurements in different epochs (e.g., NAD83(CORS96) epoch 1997.0 to NAD83(CORS96) epoch 2010.0), and the differences due to all three factors. Since most points are moving at velocities less than 0.01 mm per year in the NAD83 reference frame, epoch differences are often ignored for all but precise geodetic surveys.

Datums shifts associated with datum transformations have changed with each successive realization, as summarized in Figure 3-23, and some datums are considered functionally equivalent when combining data from different data layers, or when applying datum transformations. The WGS84(G730)

was aligned with the ITRF92 datum, so these may be substituted in datum transformations requiring no better than centimeter level accuracies. Similarly, the WGS84(G1150) and ITRF00 datums have been aligned, and may be substituted in most transformations.

While differences among the NAD83(CORSXX) and the ITRF/WGS84 datums are commonly over a meter, datum shifts internal to these groupings have become small for recent datums. Differences between NAD83(HARN) and NAD83(CORSxx) datums may be up to 20 cm, but are typically less than 4 cm, so these datum realizations may be considered equivalent if accuracy limits are above 20 cm, and perhaps as low as 4 cm. The differences among NAD83(CORS96) and NAD83(2011) are often on the order of a few centimeters, as are the differences among ITRF realizations, e.g., 91, 94, 00, 05, and 08.

There will be new datum realizations, each requiring additional transformations in the future. The ITRF datums are released every few years, requiring new transformations to existing datums each time. As of this writing, the NGS has released the NAD83(NSRS2007) datum coordinates. This is a re-analysis of state-collected points that were the basis for the NAD83(HARN) network, applying uniform, improved analysis methods. NAD83(2011) is to be released in early 2012, a nationwide adjustment of passive bench mark stations and multi-year observations at GNSS/GPS CORS stations.

Until quite recently, spatial error due to improper datum transformation has been below a detectable threshold in many analyses, so it caused few problems. GNSS receivers can now provide centimeter-level accuracy in the field, so what were once considered small discrepancies often cannot now be overlooked. As data collection accuracies improve, datum transformation errors become more apparent. The datum transformation method within any hardware or software package should be documented and the accuracy of the method known before it is

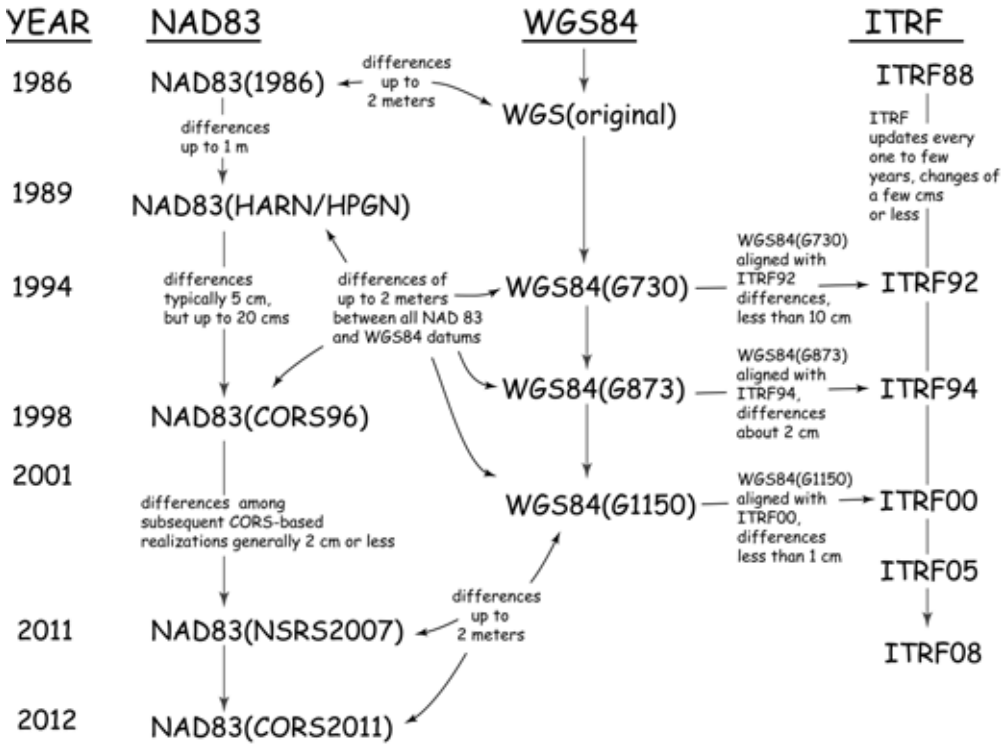


Figure 3-23: This graphic summarizes the evolution of the three main families of datums used in North America. As the datums have been adjusted, horizontal positional differences between bench mark points have varied, within the ranges shown. “Aligned” datums (e.g., WGS84(G1150) and ITRF00) may be considered equivalent for most purposes when applying datum transformations.

adopted. Unfortunately, both of these recommendations are too often ignored or only partially adopted by software vendors and users.

The NGS maintains and disseminates a list of control points in the United States (Figure 3-24), including those points used in datum definitions and adjustment. Point descriptions are provided in digital forms, including access via the world wide web (<http://www.ngs.noaa.gov>). Stations may be found based on a station name, a state and county name, a type of station (horizontal or vertical), by survey order, survey accuracy, date, or coordinate location. These stations may be used as reference points against which to check the accuracy and correctness of any data set, or as a starting point for additional surveys.

These NGS sheets may provide an estimate of the shifts associated with a datum

transformation, and so consulting one may give the specific datum shift value in any working area. For example, the values in the datum sheet in Figure 3-24 report bench mark coordinates in various datum realization that allow datum shift estimates of approximately 57.7 meters from NAD27 to NAD83(86), and 17 cm from NAD(86) to NAD83(CORS96). Similar data from a nearby station allow a calculated datum shift of approximately 1.02 m between NAD83(CORS96) and WGS84(G1150)/ITRF00. You would expect perfectly accurate data to mis-align by these amounts if the proper datum shifts were not applied.

There are a number of factors that we should keep in mind when applying datum transformations. First, changing a datum changes our best estimate of the coordinate locations of most points. These differences may be small and ignored with little penalty

in some specific instances, typically when the changes are smaller than the spatial accuracy required for our analysis. However, many datum shifts are quite large, up to tens of meters. One should know the magnitude of the datum shifts for the area and datum transformations of interest.

Second, datum transformations are estimated relationships which are developed with a specific data set and for a specific area and time. There are spatial errors in the transformations that are specific to the input and datum version. There is no generic transformation between NAD83 and WGS84. Rather, there are transformations

between specific versions of each, for example, from NAD83(96) to WGS84(1150).

Finally, GIS projects should not mix datums except under circumstances when the datum shift is small relative to the requirements of the analysis. Unless proven otherwise, all data should be converted to the same coordinate system, based on the same datum. If not, data may mis-align.

```

National Geodetic Survey, Retrieval Date = SEPTEMBER 26, 2011
OB0554 DESIGNATION - CAPE SMALL OB0554 PID - OB0554
OB0554 STATE/COUNTY- ME/SAGADAHOC USGS QUAD - PHIPPSBURG (1957)
OB0554
OB0554 *CURRENT SURVEY CONTROL
OB0554
OB0554* NAD 83(1996) - 43 46 42.87649(N) 069 50 42.26065(W) ADJUSTED
OB0554* NAVD 88 - 73. (meters) 240. (feet) SCALED
OB0554
OB0554 LAPLACE CORR- 2.33 (seconds) DEFLEC99
OB0554 GEOID HEIGHT- -25.73 (meters) GEOID03
OB0554 HORZ ORDER - FIRST
.
.
OB0554: Primary Azimuth Mark Grid Az
OB0554:SPC ME W - BURNT LEDGE JR 1866 008 26 54.5
OB0554:UTM 19 - BURNT LEDGE JR 1866 009 15 20.4
OB0554
OB0554|-----|
OB0554| PID Reference Object Distance Geod. Az |
OB0554| | | | dddmmss.s |
OB0554|
OB0555 BURNT LEDGE JR 1866 APPROX. 2.2 KM 0084015.5 | OB0554|
OB0531 MT MERRITT 2 APPROX. 1.3 KM 2083443.8 |
OB0554|-----|
OB0554
OB0554 SUPERSEDED SURVEY CONTROL
OB0554
OB0554 NAD 83(1992)- 43 46 42.87431(N) 069 50 42.25948(W) AD( ) 1
OB0554 NAD 83(1986)- 43 46 42.88001(N) 069 50 42.26497(W) AD( ) 1
OB0554 NAD 27 - 43 46 42.57400(N) 069 50 44.10300(W) AD( ) 1
.
.

```

Figure 3-24 A portion of a National Geodetic Survey control point data sheet.

Vertical Datums

Just as there are networks of well-measured points to define horizontal position, there are networks of points to define vertical position and *vertical datums*. Vertical datums are used as a reference for specifying heights. Much like horizontal datums, they are established through a set of painstakingly surveyed control points. These point elevations are precisely measured, initially through a set of optical surface measurements, but more recently using GPS, laser, satellite, and other measurement systems. Establishing vertical datums also requires estimating the strength and direction of the gravitational force near the surface of the Earth.

In its simplest definition, a vertical datum is a reference that we use for measuring heights. As noted in the geoid section on page 78, we use a geoid as a reference surface, and specify the orthometric heights as the elevations of points on the Earth's surface above the geoid. We first establish a specific geoid through a set of gravity measurements and then augment this with precise vertical height measurements at points across the globe to establish a set of vertical bench marks, against which we can conveniently measure all other heights. The vertical datum is the set of points, with heights, relative to a specific geoid.

Leveling surveys are among the oldest methods for establishing a vertical point. Distances and elevation differences are precisely measured from an initial point to other points, establishing height differentials. Early leveling surveys were performed with the simplest of instruments, including a plumb bob to establish leveling posts, and a simple liquid level to establish horizontal lines. Early surveys used an approach known as *spirit leveling*. Horizontal rods were placed between succeeding leveling posts across the landscape to physically measure height differences (Figure 3-25).



Figure 3-25: Early surveys used level bars placed on vertical posts, simple but effective technology.

The number, accuracy, and extent of leveling surveys increased substantially in the 18th and 19th centuries. Epic surveys that lasted decades were commissioned, such as the Great Arc, from southern India to the Himalayas. These surveys were performed at substantial capital and human expense; in one portion of the Great Arc more than 60% of the field crews died due to illness and mishaps over a six year period. Surface leveling provided most height measurements for vertical datums until the mid to late 20th century, when a variety of satellite-based methods were introduced.

Most leveling surveys from the late 1700s through the mid 20th century employed *trigonometric leveling*. This method uses optical instruments and trigonometry to measure changes in height, as shown in Figure 3-26. Surface distance along the slope was measured to avoid the tedious process of establishing vertical posts and leveling rods. The vertical angle was also measured from a known station to an unknown station. The angle was typically measured with a small telescope fitted with a precisely scribed angle gauge. The gauge could be referenced to zero at a horizontal position, usually with an integrated bubble level, or later, with an electronic level. Surface distance would then be combined with

the measured vertical angle to calculate the horizontal and vertical distances. Early surveys measured surface distance along the slope with ropes, metal chains, and steel tapes, but these physical devices have largely been replaced by improved optical methods, or by laser-based methods.

Early national leveling surveys used the concept of mean sea level as a zero, or base height. Sea-level heights vary over time, mostly due to tides, but also due to changes in weather systems, currents, temperature and salinity. Mean sea level at a gauge may be calculated after a sufficient period of time, typically over at least the 19-year cycle of tidal variation. Monumented points were established on rocks, docks, or other ocean-side fixed objects near the gauges, and the height of these starting points could then be measured via leveling to the nearby ocean tidal stations. Precise leveling was then extended landward from these oceanside points to measure heights cross-country. All leveled heights could then be tied to a mean sea level through this vertical measurement network.

Note that we said “a” mean sea level, because mean sea level isn’t the same everywhere. Mean sea level, even averaged over

several decades, varies across the globe due to several factors, for example, persistent differences in water density with temperature and salinity, or regular ocean currents, which may persistently raise or lower the surface in ocean regions. This means the mean sea level is not constant relative to the geoid or ellipsoid, and will be different at Miami than New York. Modern vertical datums do not use mean sea level across many stations as a reference in part because of this variation in mean sea level across the Earth. While most people describe mountain summit elevations or other heights as above mean sea level, geodesists and GIS professionals do not. We use a set of precisely surveyed base bench marks with orthometric heights referenced to the Earth’s geoid.

As with horizontal datums, the primary vertical datums in use have changed through time as the number, distribution, and accuracy of vertical survey points have increased. Geodetic leveling surveys began in the U.S. in the 1850s, initially focusing on the East Coast and Great Lakes region, and extended across the U.S. between 1877 and 1900. Periodic adjustments harmonized measurements, identified and removed large errors, and distributed small discrepancies

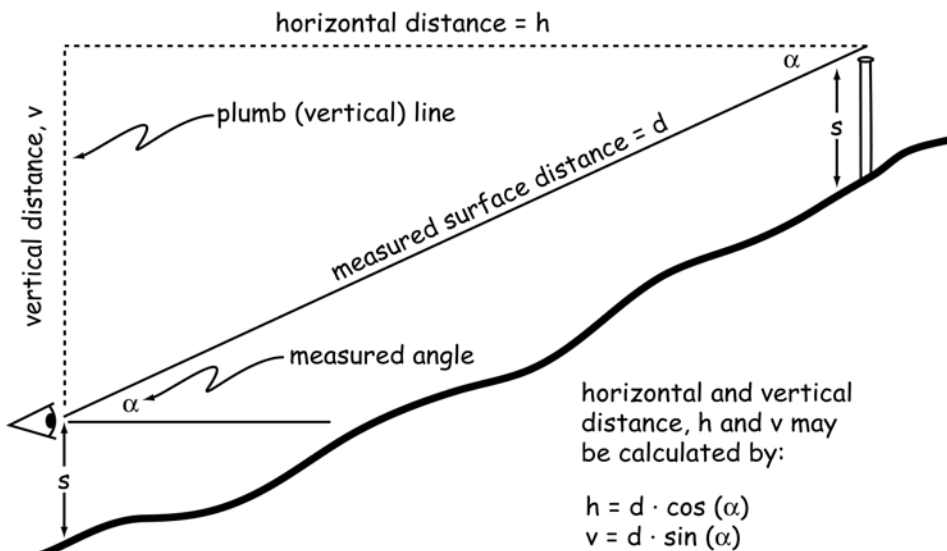


Figure 3-26: Leveling surveys often employ optical measurements of vertical angle (α) with measurements of surface distance (d) and knowledge of trigonometric relationships to calculate horizontal distance (h) and vertical distance (v).

```

National Geodetic Survey, Retrieval Date = FEBRUARY 12, 2011
FB2737 *****
FB2737 CBN - This is a Cooperative Base Network Control Station.
FB2737 DESIGNATION - MITCHELL 2
FB2737 PID - FB2737
FB2737 STATE/COUNTY- NC/YANCEY
FB2737 USGS QUAD - MT MITCHELL (1946)
FB2737
FB2737 *CURRENT SURVEY CONTROL
FB2737
FB2737 *-----*
FB2737* NAD 83(NSRS2007)- 35 45 53.76569(N) 082 15 54.34377(W) ADJUSTED
FB2737* NAVD 88 - 2048.27 (meters) 6720. (feet) LEVELING
FB2737 *-----*
FB2737 X - 697,569.032 (meters) COMP
FB2737 Y - -5,135,766.495 (meters) COMP
FB2737 Z - 3,708,239.126 (meters) COMP
FB2737 LAPLACE CORR- -6.76 (seconds) DEFLEC99
FB2737 ELLIP HEIGHT- 2018.409 (meters) (02/10/07) GPS OBS
FB2737 GEOID HEIGHT- -30.00 (meters) GEOID03
    
```

Figure 3-27: A portion of a data sheet for a vertical control bench mark.

among stations. These vertical adjustments were conducted in 1899, 1903, 1907, and 1912, relating all measured heights to between five and nine precisely measured tidal gauges.

The first continental vertical datum in North America was the *National Geodetic Vertical Datum* of 1929, also referred to as NGVD29. Vertical leveling was adjusted to 26 tidal gauges, including 5 in Canada, based on local mean sea level at each of the gauges. Geodesists realized that mean sea level varied across the continent, but assumed these differences would be similar or smaller than measurement errors. They wanted to avoid confusion caused by seaside bench marks having heights that differed from mean sea level.

Vertical measurements continued from the 1920s through the 1980s, resulting in the *North American Vertical Datum of 1988* (NAVD88), and many monumented control points have vertical heights reported in NAVD88 (Figure 3-27). The 1988 datum is based on over 600,000 kilometers (360,000 miles) of control leveling performed since 1929, and also reflects geologic crustal movements or subsidence that may have changed bench mark elevation. NAVD88 was fixed relative to only one tidal station, in the town of Rimouski, Quebec, because

improved methods meant measurement errors were much smaller than differences in mean sea level among stations. Surface heights in this datum are not based on mean sea level, because to do so across the set of sea-side benchmarks would require additional warping that would degrade the measurements.

Improved geoid models have been developed concurrently with these newer vertical datums. For example, at this writing, the most current model, GEOID03, integrated nearly 15,000 vertical bench marks to estimate geoidal and orthometric heights. These heights are reported on NGS data sheets for vertical bench marks (Figure 3-27), noting the vertical datum (here NAVD88), the geoid model (GEOID03), the orthometric height (here 2048.27 meters), and the ellipsoidal and geoidal heights.

Dynamic Heights

We must discuss one final kind of height, called a dynamic height, because they are important for certain applications, and are often listed on NGS data sheets and elsewhere. Dynamic heights measure the change in gravitational pull from a given equipotential surface. Dynamic heights are important when interested in water levels

and flows across elevations. Points that have the same dynamic heights can be thought of as being at the same water level. Perhaps a bit surprisingly, points with the same dynamic heights often have different orthometric heights (Figure 3-28). To be clear, two distinct points at water's edge on a large lake often do not have the same elevations, that is, they are different orthometric heights above our reference geoid. Since orthometric heights are our bench mark for specifying elevation, water may flow from one point to another, even though those points have the same elevation.

To understand why water may flow between points with the same elevation (orthometric heights), it is important to remember how orthometric heights are defined. An orthometric height is the distance, in the direction of gravitational pull, from the geoid up to a point. But remember, the geoid is a specified gravity value, an "equipotential" surface, where the pull of gravity is at some specified level. As we move up from the geoid toward the surface, we pass through other equipotential sur-

faces, each at a slightly weaker gravitational pull or force, until we arrive at the surface point.

There are two key observations here. First, water spreads out to level across an equipotential surface, absent wind, waves, and other factors. The water level in a still bathtub, pond, or lake is at the same equipotential surface at one end as another. Gravity pulls down on the surface to ensure it conforms to an equipotential surface. Second, the equipotential surfaces are closer together when nearer the mass center of Earth. The equipotential surfaces converge, or become "denser" the closer you are to the center of the Earth.

Because of these two facts, and because the Earth's polar radius is less than the equatorial radius, the orthometric heights of the water surface on large lakes are usually different at the north and south ends. For example, as you move further north in the northern hemisphere, the equipotential surfaces converge due to a decreasing distance from the mass center of the Earth, and the pull of gravity increase (Figure 3-28). An

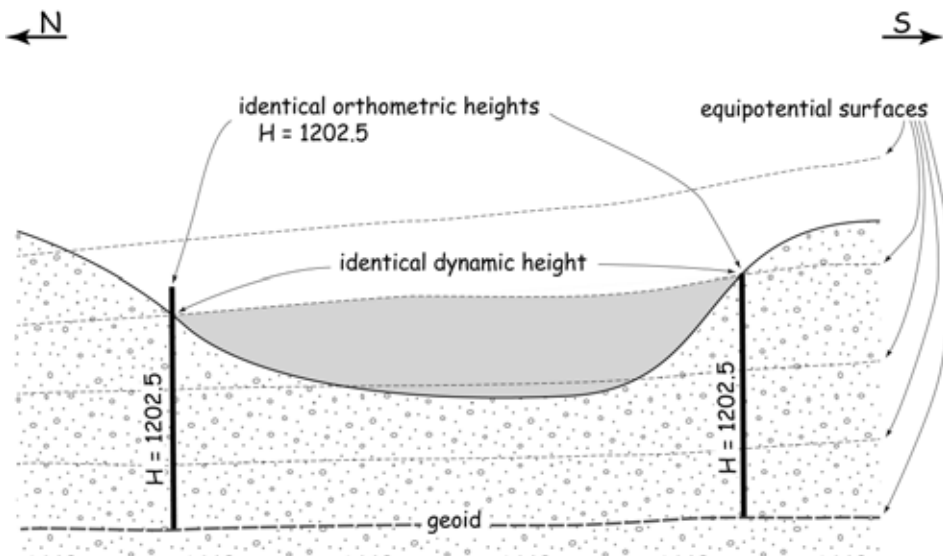


Figure 3-28: An illustration of how dynamic heights and orthometric heights may differ, and how equal orthometric heights may correspond to different heights above the water level on a large lake. Because equipotential surfaces converge, the orthometric height at the water level at the northern extreme of a lake will have different orthometric heights. Dynamic heights and water levels are equal across an equipotential surface.

orthometric height is a fixed height above the geoidal surface, so the northern orthometric height will pass through more equipotential surfaces than the same orthometric height at a more southerly location. Water follows an equipotential surface, so an orthometric height of the water level at the south end of the lake will be higher than at the north end. For example, in the Great Lakes of North America, the orthometric height corresponding to water level at the south end of Lake Michigan is approximately 15 cm higher than the water level at the north end.

Dynamic heights are most often used when we're interested in relative heights for water levels, particularly over large lakes or connected water bodies. Because equal dynamic heights are at the same water level, we can use them when interested in accurately representing hydrologic drop, head, pressure, and other variables related to water levels across distances. But these differences should be confusing when observing benchmark or sea level heights, and underscore that our height reference is not mean sea level, but rather an estimated geoidal surface.

Control Accuracy Specification

In most cases the horizontal datum control points are too sparse to be sufficient for all needs in GIS data development. For example, precise point locations may be required when setting up a GNSS receiving station, to georegister a scanned photograph or other imagery, or as the basis for a detailed subdivision or highway survey. It is unlikely there will be more than one or two datum points within any given work area. Because a denser network of known points is required for many projects, datum points are often used as a starting locations for additional surveying. These smaller area surveys increase the density of precisely known points. The quality of the point locations

depends on the quality of the intervening survey.

The Federal Geodetic Control Committee of the United States (FGCC) has published a detailed set of survey accuracy specifications. These specifications set a minimum acceptable accuracy for surveys and establish procedures and protocols to ensure that the advertised accuracy has been obtained. The FGCC specifications establish a hierarchy of accuracy. First order survey measurements are accurate to within 1 part in 100,000. This means the error of the survey is no larger than one unit of measure for each 100,000 units of distance surveyed. The maximum horizontal measurement error of a 5,000 meter baseline (about 3 miles) would be no larger than 5 centimeters (about 2 inches). Accuracies are specified by Class and Orders, down to a Class III, 2nd order point with an error of no more than 1 part in 5,000.

Map Projections and Coordinate Systems

Datums tell us the latitudes and longitudes of a set of points on an ellipsoid. We need to transfer the locations of features measured with reference to these datum points from the curved ellipsoid to a flat map. A *map projection* is a systematic rendering of locations from the curved Earth surface onto a flat map surface. Points are “projected” from the Earth surface and onto the map surface.

Most map projections may be viewed as sending rays of light from a projection source (Figure 3-29). Rays radiate from a source to intersect both the ellipsoid surface and the map surface. The rays specify where each point from the ellipsoid surface is placed on the map surface. In some projections the source is not a single point; however the basic process involves the systematic transfer of points from the curved ellipsoidal surface to a flat map surface.

Distortions are unavoidable when making flat maps, because as we’ve said, locations are projected from a complexly curved Earth surface to a flat or simply curved map surface. Portions of the rendered Earth sur-

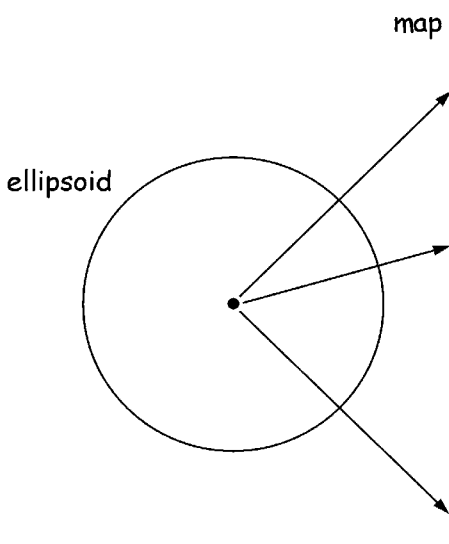


Figure 3-29: A conceptual view of a map projection.

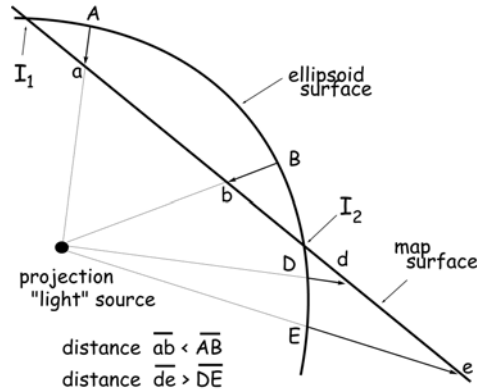


Figure 3-30: Distortion during map projection.

face must be compressed or stretched to fit onto the map. This is illustrated in Figure 3-30, a side view of a projection from an ellipsoid onto a plane. The map surface intersects the Earth at two locations, I_1 and I_2 . Points toward the edge of the map surface, such as D and E, are stretched apart. The scaled map distance between d and e is greater than the distance from D to E measured on the surface of the Earth. More simply put, the distance along the map plane is greater than the corresponding distance along the curved Earth surface. Conversely, points such as A and B that lie in between I_1 and I_2 would appear compressed together. The scaled map distance from a to b would be less than the surface-measured distance from A to B. Distortions at I_1 and I_2 are zero.

Figure 3-30 demonstrates a few important facts. First, distortion may take different forms in different portions of the map. In one portion of the map features may be compressed and exhibit reduced areas or distances relative to the Earth’s surface measurements, while in another portion of the map areas or distances may be expanded. Second, there are often a few points or lines where distortions are zero and where length, direction, or some other geometric property is preserved. Finally, distortion is usually less near the points or lines of intersection,

where the map surface intersects the imaginary globe. Distortion usually increases with increasing distance from the intersection points or lines.

Different map projections may distort the globe in different ways. The projection source, represented by the point at the middle of the circle in Figure 3-30, may change locations. The surface onto which we are projecting may change in shape, and we may place the projection surface at different locations at or near the globe. If we change any of these three factors, we will change how or where our map is distorted. The type and amount of projection distortion may guide selection of the appropriate projection or limit the area projected.

Figure 3-31 shows an example of distortion with a projection onto a planar surface. This planar surface intersects the globe at a line of true scale, the solid line labeled as the standard circle shown in Figure 3-31. Distortion increases away from the line of true

scale, with features inside the circle compressed or reduced in size, for a negative scale distortion. Conversely, features outside the standard circle are expanded, for a positive scale distortion. Calculations show a scale error of -1% near the center of the circle, and increasing scale error in concentric bands outside the circle to over 2% near the outer edges of the projected area.

An approximation of the distance distortion may be obtained for any projection by comparing grid coordinate distances to *great circle distances*. A great circle distance is a distance measured on the ellipsoid and in a plane through the Earth's center. This planar surface intersects the two points on the Earth's surface and also splits the spheroid into two equal halves (Figure 3-32). The smallest great circle distance is the shortest path between two points on the surface of the ellipsoid, and by approximation, Earth.

As noted earlier, a straight line between two points on the projected map is likely not

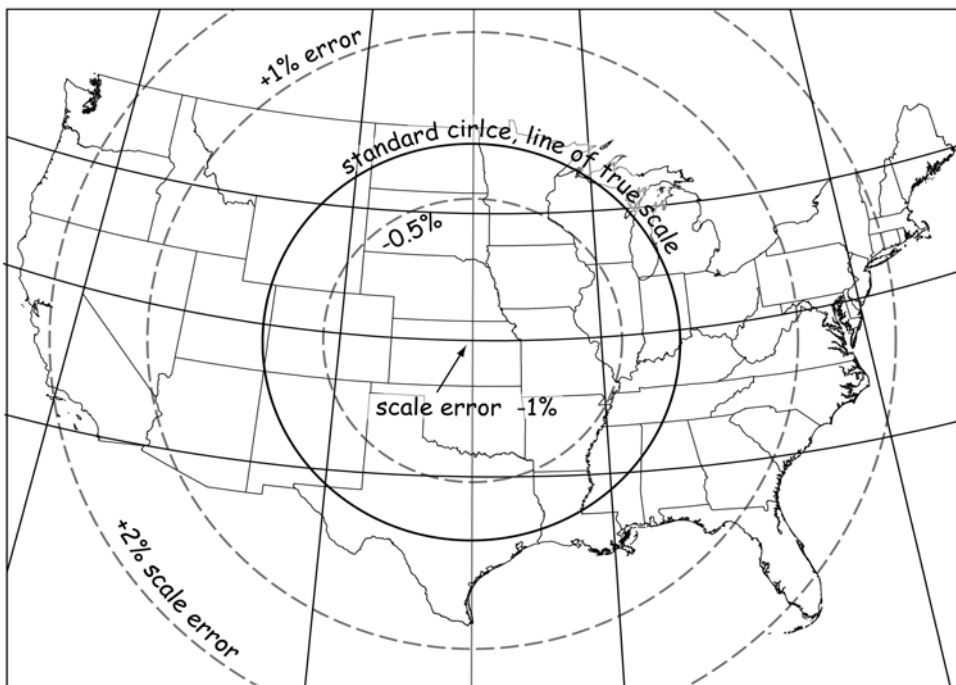
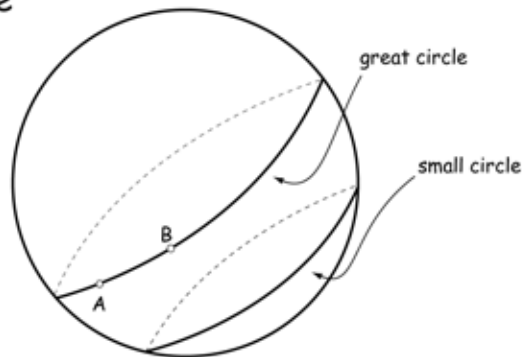


Figure 3-31: Approximate error due to projection distortion for a specific oblique stereographic projection. A plane intersects the globe at a standard circle. This standard circle defines a line of true scale, where there is no distance distortion. Distortion increases away from this line, and varies from -1% to over 2% in this example (adapted from Snyder, 1987).

Great Circle Distance

Consider two points on the Earth's surface, A with geographic coordinates (lat., lon.) (ϕ_A, λ_A) , and

B, with geographic coordinates (ϕ_B, λ_B)



The great circle distance from point A to point B is given by the formula:

$$d = r \cdot \cos^{-1}[(\cos(\phi_A)\cos(\phi_B)\cos(\lambda_A - \lambda_B) + \sin(\phi_A)\sin(\phi_B))],$$

where d is the shortest distance on the surface of the Earth from A to B, and r is the Earth's radius, approximately 6378 km.

This formula may be used to find the distance distortion caused by a projection between two points, for example, between Ursine and Moab, Utah, when using UTM Zone 12N coordinates, NAD83?

Great circle distance:

Latitude, longitude of Ursine, Utah = $37.98481^\circ, -114.216944^\circ$

Latitude, longitude of Moab, Utah = $38.57361^\circ, -109.551111^\circ$

$$\begin{aligned} d &= 6378 \cdot \cos^{-1}[(\cos(37.98481)\cos(38.57361)\cos(-114.216944 - 109.551111) + \\ &\quad \sin(37.98481)\sin(38.57361))] \\ &= 412.906 \text{ km} \end{aligned}$$

Grid distance (UTM Zone 12N coordinates):

Grid coordinates of Ursine, Utah = 217,529.8, 4,208,972.8

Grid coordinates of Moab, Utah = 626,239.2, 4,270,405.9

$$\begin{aligned} dg &= [(X_A - X_B)^2 + (Y_A - Y_B)^2]^{0.5} \\ &= [(217,529.8 - 626,239.2)^2 + (4,208,972.8 - 4,270,405.9)^2]^{0.5} \\ &= 413.300 \text{ km} \end{aligned}$$

distortion is $412.906 - 413.300 = -0.394$ km, or a 394 meter lengthening

Figure 3-32: Example calculation of the distance distortion due to a map projection. The great circle and grid distances are compared for two points on the Earth's surface, the first measuring along the curved surface, the second on the projected surface. The difference in these two measures is the distance distortion due to the map projection. Calculations of the great circle distances are approximate, due to the assumption of a spheroidal rather than ellipsoidal Earth, but are very close.

to be a straight line on the surface of the Earth, and is not the shortest distance between two points when traveling on the surface of the Earth. Conversely, the shortest distance between points when traveling on the surface of the Earth is likely to appear as a curved line on a projected map. The distortion is imperceptible for large scale maps and over short distances, but exists for most lines.

Figure 3-33 illustrates straight line distortion. This figure shows the shortest distance path between Adelaide, Australia, and Tokyo, Japan. Tokyo lies almost due north of Adelaide, and the shortest path approximates a line of longitude, by definition a great circle path. This shortest path is distorted and appears curved by the projection used for this map.

The magnitude of this distortion may be approximated by simple formulas (Figure 3-32). Coordinates may be identified for any

two points in the grid system, and the Pythagorean formula used to calculate distance between the two points. The resulting distance will be expressed in the grid coordinate system, and therefore will include the projection distortion. The distance may also be calculated for a great circle route along the spheroid surface. This calculation will approximate the unprojected distance, measured on the surface of the Earth. This is only an approximation, as we know from the previous section, because the Earth is shaped more like an ellipsoid, and has geoidal undulations. However, the approximation is quite accurate, generally off by less than a few parts per tens of thousands over several hundred kilometers. The great circle and grid coordinate distance may then be compared to estimate the distance distortion (Figure 3-32).

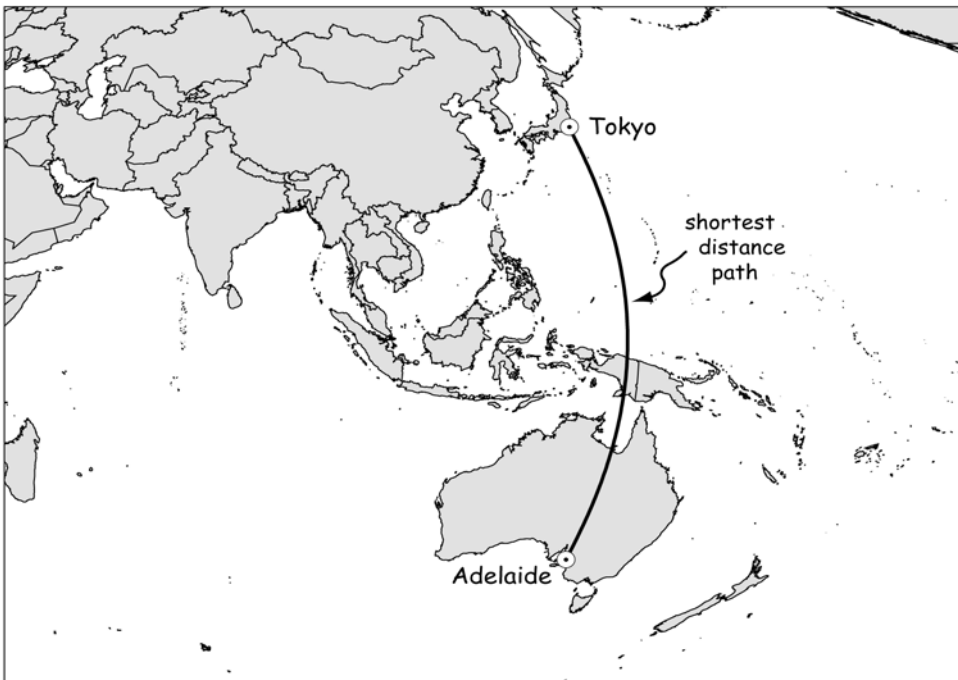


Figure 3-33: Curved representations of straight lines are a manifestation of projection distortion. A great circle path, shown above, is the shortest route when traveling on the surface of the Earth. This path appears curved when plotted on this sinusoidal projection.

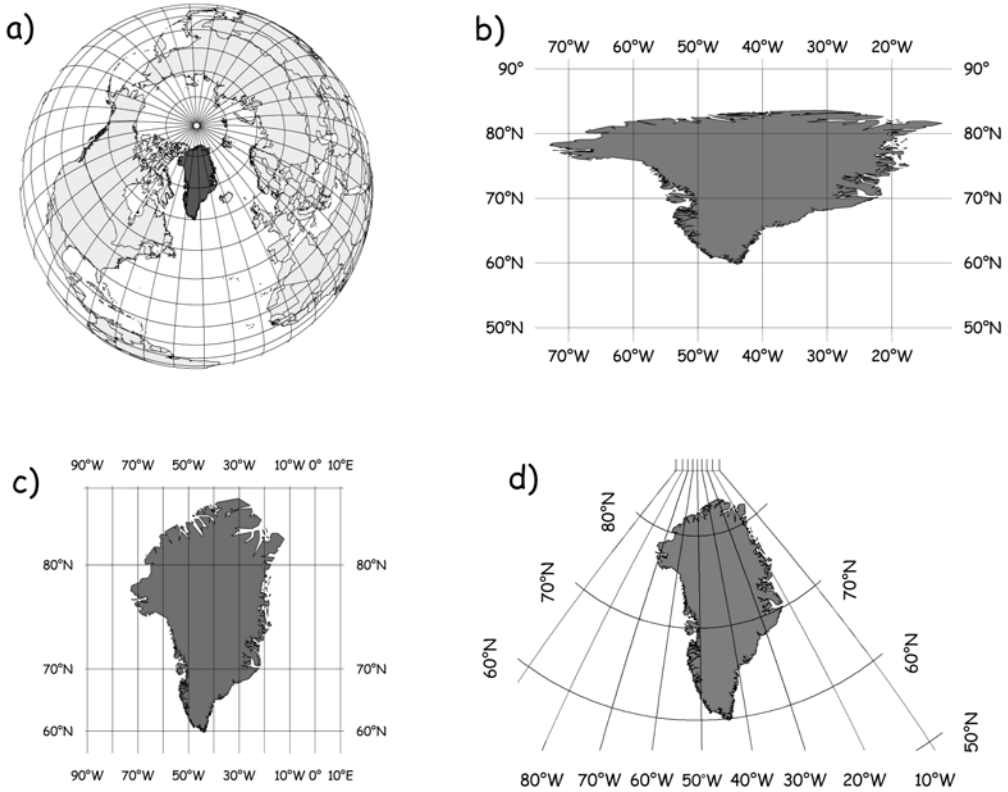


Figure 3-34: Map projections can distort the shape and area of features, as illustrated with these various projections of Greenland, from a) approximately unprojected, b) geographic coordinates on a plane, c) a Mercator projection, and d) a transverse Mercator projection.

Projections may also substantially distort the shape and area of polygons. Figure 3-34 shows various projections for Greenland, from an approximately “unprojected” view from space through geographic coordinates cast on a plane, to Mercator and transverse Mercator projections. Note the changes in size and shape of the polygon depicting Greenland.

Most map projections are based on a *developable surface*, a geometric shape onto which the Earth surface locations are projected. Cones, cylinders, and planes are the most common types of developable surfaces. A plane is already flat, and cones and cylinders may be mathematically “cut” and “unrolled” to develop a flat surface (Figure 3-35). Projections may be characterized according to the developable surface, for

example, as *conic* (cone), *cylindrical* (cylinder), and *azimuthal* (plane). The orientation of the developable surface may also change among projections, for example, the axis of a cylinder may coincide with the poles (equatorial) or the axis may pass through the equator (transverse).

Note that while the most common map projections used for spatial data in a GIS are based on a developable surface, many map projections are not. Projections with names such as pseudocylindrical, Mollweide, sinusoidal, and Goode homolosine are examples. These projections often specify a direct mathematical projection from an ellipsoid onto a flat surface. They use mathematical forms not related to cones, cylinders, planes, or other three-dimensional figures, and may change the projection surface for different

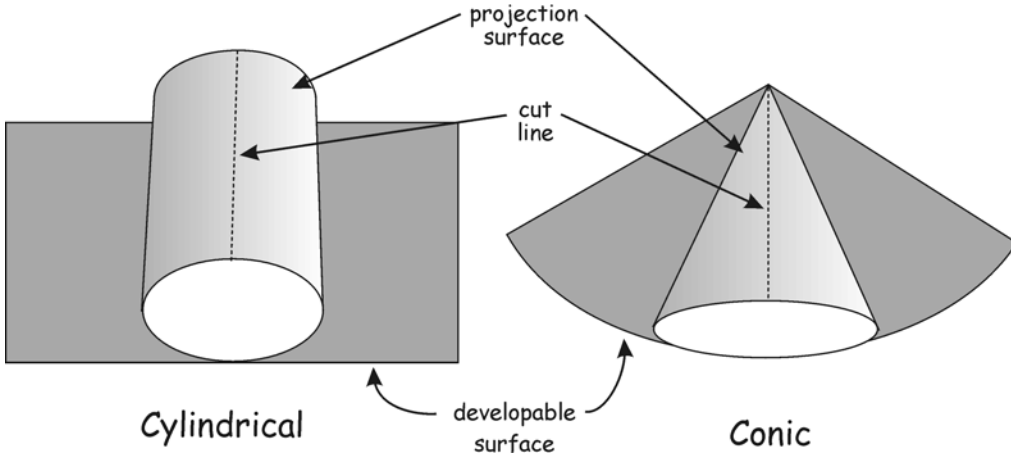


Figure 3-35: Projection surfaces are derived from curved “developable” surfaces that may be mathematically “unrolled” to a flat surface.

parts of the globe. For example, projections such as the Goode homolosine projection are formed by fusing two or more projections along specified line segments. These projections use complex rules and breaks to reduce distortion for many continents.

We typically have to specify several characteristics when we specify a map projection. For example, for an azimuthal projection we must specify the location of the projection center (Figure 3-36) and the location and orientation of the plane onto which the globe is projected. Azimuthal projections are often tangent to (just touch) the ellipsoid at one point, and we must specify the location of this point. A projection center (“light” source location) must also be specified, most often placed at one of three locations. The projection center may be at the center of the ellipsoid (a *gnomonic* projection), at the antipodal surface of the ellipsoid (diametrically opposite the tangent point, a *stereographic* projection), or at infinity (an *orthographic* projection). Scale factors, the location of the origin of the coordinate system, and other projection parameters may be required. Defining characteristics must be specified for all projections, such as the size and orientation of a cone in a conic projection, or the size, intersection properties, and orientation of a cylinder in a cylindrical projection.

Note that the use of a projection defines a projected coordinate system and hence typically adds a third version of North to our description of geography. We have already described magnetic north, towards which a compass points, and geographic north, the pole around which the globe revolves (Figure 3-12). We must add *grid north* to these, defined as the direction of the Y axis in the projection. Grid north is often defined by some meridian in the projection, often known as the central meridian. Grid north is typically not the same as geographic or magnetic north over most of the projected area.

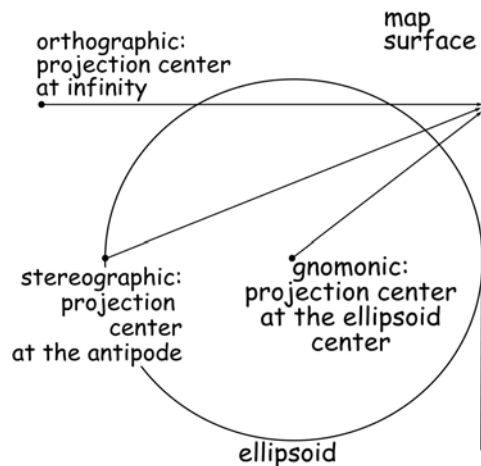


Figure 3-36: The projection center of map projections is most often placed at the center of the ellipsoid, or the antipode, or at infinity.

Common Map Projections in GIS

There are hundreds of map projections used throughout the world, however most spatial data in GIS are specified using a relatively small number of projection types.

The Lambert conformal conic and the transverse Mercator are among the most common projection types used for spatial data in North America, and much of the world (Figure 3-37). Standard sets of projections have been established from these two basic types. The Lambert conformal conic (LCC) projection may be conceptualized as a cone intersecting the surface of the Earth, with points on the Earth's surface projected onto the cone. The cone in the Lambert conformal conic intersects the ellipsoid along two arcs, typically parallels of latitude, as shown in Figure 3-37 (top left). These lines of intersection are known as *standard parallels*.

Distortion in a Lambert conformal conic projection is typically smallest near the standard parallels, where the developable surface intersects the Earth. Distortion increases in a complex fashion as distance from these intersection lines increases. This characteristic is illustrated at the top right and bottom of Figure 3-37. Circles of a constant 5 degree radius are drawn on the projected surface at the top right, and approximate lines of constant distortion and a line of true scale are shown in Figure 3-37, bottom. Distortion decreases towards the standard parallels, and increases away from these lines. Those farther away tend to be more distorted. Distortions can be quite severe, as illustrated by the apparent expansion of southern South America.

Note that sets of circles in an east-west row are distorted in the Lambert conformal conic projection (Figure 3-37, top right). Those circles that fall between the standard parallels exhibit a uniformly lower distortion than those in other portions of the projected map. One property of the Lambert conformal conic projection is a low-distortion band running in an east-west direction between the standard parallels. Thus, the Lambert

conformal conic projection is often used for areas that are larger in an east-west than a north-south direction, as there is little added distortion when extending the mapped area in the east-west direction.

Distortion is controlled by the placement and spacing of the standard parallels, the lines where the cone intersects the globe. The example in Figure 3-37 shows parallels placed such that there is a maximum distortion of approximately 1% midway between the standard parallels. We reduce this distortion by moving the parallels closer together, but at the expense of reducing the area mapped at this lower distortion level.

The transverse Mercator is another common map projection. This map projection may be conceptualized as enveloping the Earth in a horizontal cylinder, and projecting the Earth's surface onto the cylinder (Figure 3-38). The cylinder in the transverse Mercator commonly intersects the Earth ellipsoid along a single north-south tangent, or along two *secant* lines, noted as the lines of true scale in Figure 3-38. A line parallel to and midway between the secants is often called the central meridian. The central meridian extends north and south through transverse Mercator projections.

As with the Lambert conformal conic, the transverse Mercator projection has a band of low distortion, but this band runs in a north-south direction. Distortion is least near the line(s) of intersection. The graph at the top right of Figure 3-38 shows a transverse Mercator projection with the central meridian (line of intersection) at 0 degrees longitude, traversing western Africa, eastern Spain, and England. Distortion increases markedly with distance east or west away from the intersection line, for example, the shape of South America is severely distorted in the top right of Figure 3-38. The drawing at the bottom of Figure 3-38 shows lines estimating approximately equal scale distortion for a transverse Mercator projection centered on the USA. Notice that the distortion increases as distance from the two lines of intersection increases. Scale distortion error may be maintained below any thresh-

Lambert Conformal Conic

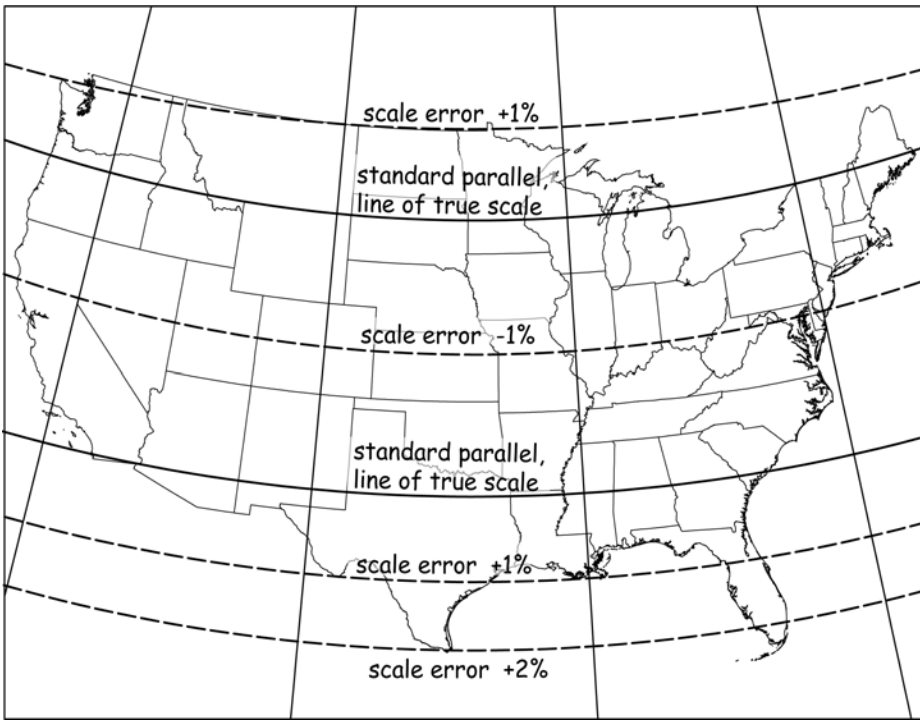
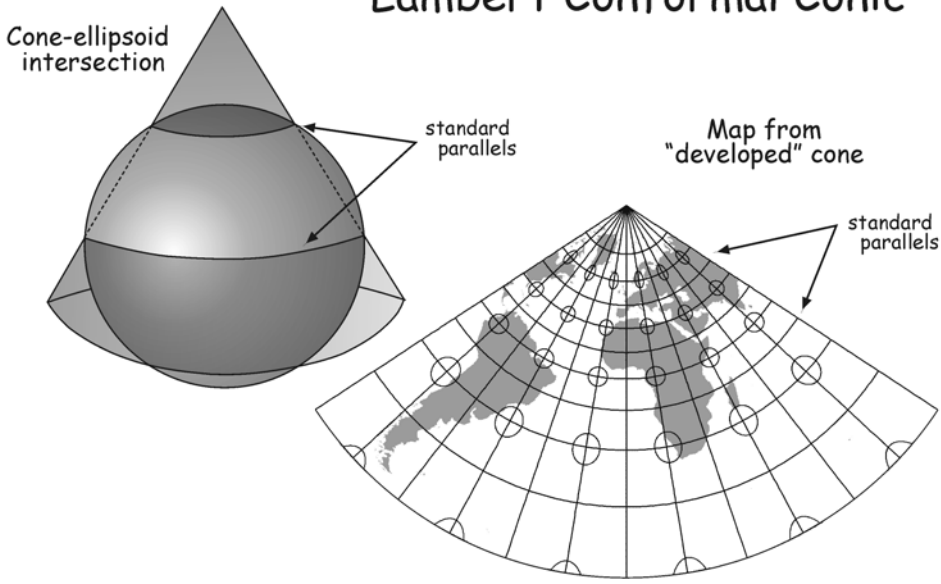


Figure 3-37: Lambert conformal conic (LCC) projection (top) and an illustration of the scale distortion associated with the projection. The LCC is derived from a cone intersecting the ellipsoid along two standard parallels (top left). The “developed” map surface is mathematically unrolled from the cone (top right). Distortion is primarily in the north-south direction, and is illustrated in the developed surfaces by the deformation of the 5-degree diameter geographic circles (top) and by the lines of approximately equal distortion (bottom). Note that there is no scale distortion where the standard parallels intersect the globe, at the lines of true scale (bottom, adapted from Snyder, 1987).

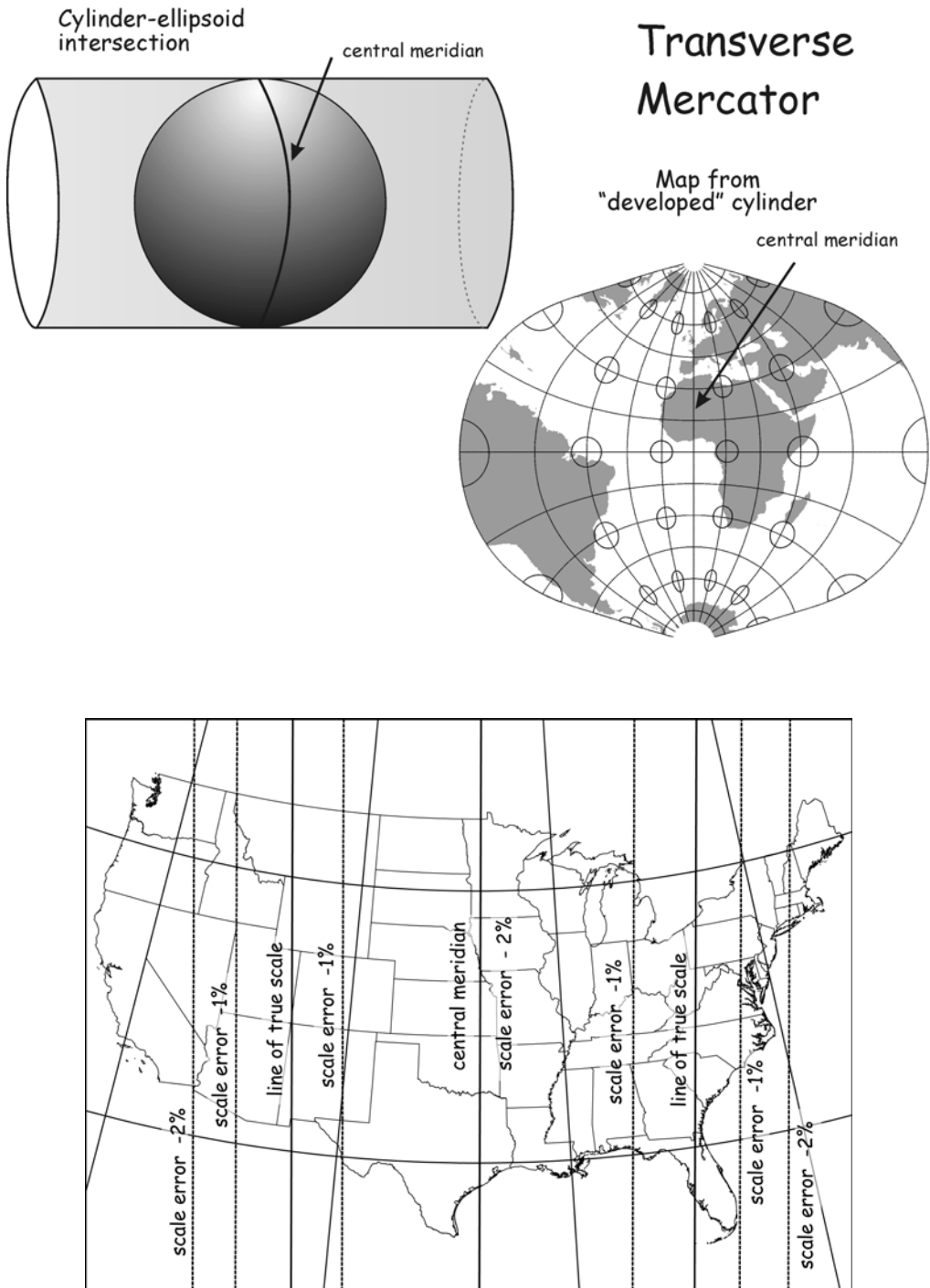


Figure 3-38: Transverse Mercator (TM) projection (top), and an illustration of the scale distortion associated with the projection (bottom). The TM projection distorts distances in an east-west direction, but has relatively little distortion in a north-south direction. This TM intersects the sphere along two lines, and distortion increases with distance from these lines (bottom, adapted from Snyder, 1987).

old by ensuring the mapped area is close to these two secant lines intersecting the globe. Transverse Mercator projections are often used for areas that extend in a north-south direction, as there is little added distortion extending in that direction.

Different projection parameters may be used to specify an appropriate coordinate system for a region of interest. Specific standard parallels or central meridians are chosen to minimize distortion over a mapping area. An origin location, measurement units, x and y (or northing and easting) offsets, a scale factor, and other parameters may also be required to define a specific projection. Once a projection is defined, the coordinates of every point on the surface of the Earth may be determined, usually by a closed-form or approximate mathematical formula.

The State Plane Coordinate System

The State Plane Coordinate System is a standard set of projections for the United States. The State Plane coordinate system specifies positions in Cartesian coordinate systems for each state. There are one or more zones in each state, with slightly different projections in each State Plane zone

(Figure 3-39). Multiple State Plane zones are used to limit distortion errors due to map projections.

State Plane systems greatly facilitate surveying, mapping, and spatial data development in a GIS, particularly when whole county or larger areas are involved. The State Plane system provides a common coordinate reference for horizontal coordinates over county to multi-county areas while limiting distortion error to specified maximum values. Most states have adopted zones such that projection distortions are kept below one part in 10,000. Some states allow larger distortions (e.g., Montana, Nebraska). State Plane coordinate systems are used in many types of work, including property surveys, property subdivisions, large-scale construction projects, and photogrammetric mapping, and the zones and state plane coordinate system are often adopted for GIS.

One State Plane projection zone may suffice for small states. Larger states commonly require several zones, each with a different projection, for each of several geographic zones of the state. For example Delaware has one State Plane coordinate zone, while California has 6, and Alaska has 10 State Plane coordinate zones, each corresponding to a different projection within the

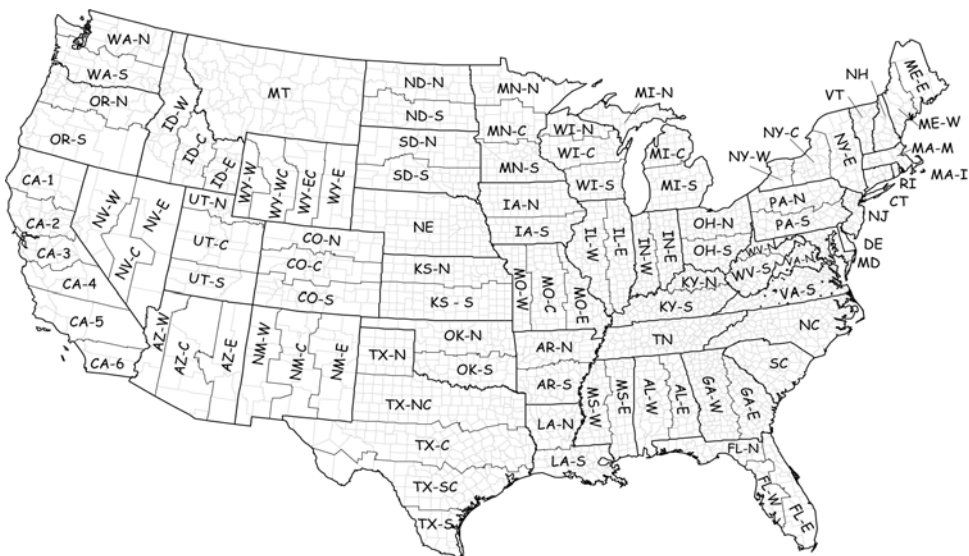


Figure 3-39: State plane zone boundaries, NAD83.

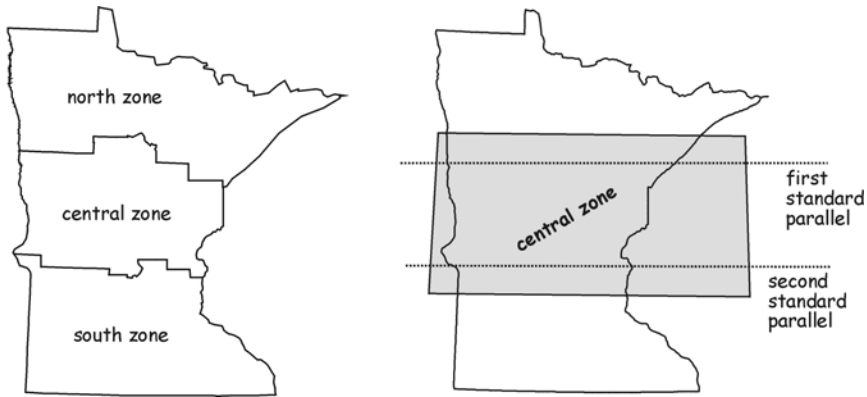


Figure 3-40: The State Plane zones of Minnesota, and details of the standard parallel placement for the Minnesota central State Plane zone.

state. Zones are added to a state to ensure acceptable projection distortion within all zones (Figure 3-40, left). Zone boundaries are defined by county, parish, or other municipal boundaries. For example, the Minnesota south/central zone boundary runs approximately east-west through the state along defined county boundaries (Figure 3-40, left).

The State Plane coordinate system is based on two types of map projections: the Lambert conformal conic and the transverse Mercator projections. Because distortion in a transverse Mercator increases with distance from the central meridian, this projection type is most often used with states that have a long north-south axis (e.g., Illinois or New Hampshire). Conversely, a Lambert conformal conic projection is most often used when the long axis of a state is in the east-west direction (e.g. North Carolina and Virginia). When computing the State Plane coordinates, points are projected from their geodetic latitudes and longitudes to x and y coordinates in the State Plane systems.

The Lambert conformal conic projection is specified in part by two standard parallels that run in an east-west direction. A different set of standard parallels is defined for each State Plane zone. These parallels are placed at one-sixth of the zone width from the north and south limits of the zone (Figure 3-40, right). The zone projection is defined by

specifying the standard parallels and a central meridian that has a longitude near the center of the zone. This central meridian points in the direction of geographic north, however all other meridians converge to this central meridian, so they do not point to geographic north. The Lambert conformal conic is used to specify projections for State Plane zones for 31 states.

As noted earlier, the transverse Mercator specifies a central meridian. This central meridian defines grid north in the projection. A line along the central meridian points to geographic north, and specifies the Cartesian grid direction for the map projection. All parallels of latitude and all meridians except the central meridian are curved for a transverse Mercator projection, and hence these lines do not parallel the grid x or y directions. The transverse Mercator is used for 22 State Plane systems (the sum of states is greater than 50 because both the transverse Mercator and Lambert conformal conic are used in some states, e.g., Florida).

Finally, note that more than one version of the State Plane coordinate system has been defined. Changes were introduced with the adoption of the North American Datum of 1983. Prior to 1983, the State Plane projections were based on NAD27. Changes were minor in some cases, and major in others, depending on the state and State Plane zone. Some states, such as South Carolina,

Nebraska, and California, dropped zones between the NAD27 and NAD83 versions (Figure 3-41). Others maintained the same number of State Plane zones, but changed the projection by the placement of the meridians, or by switching to a metric coordinate system rather than one using feet, or by shifting the projection origin. State Plane zones are sometimes identified by the Federal Information Processing System (FIPS) codes, and most codes are similar across NAD27 and NAD83 versions. Care must be taken when using older data to identify the version of the State Plane coordinate system used because the FIPS and State Plane zone designators may be the same, but the projection parameters may have changed from NAD27 to NAD83.

Conversion among State Plane projections may be additionally confused by the various definitions used to translate from feet to meters. The metric system was first developed during the French Revolution in the late 1700s, and it was adopted as the official unit of distance in the United States, by

the initiative of Thomas Jefferson. President Jefferson was a proponent of the metric system because it improved scientific measurements, was based on well-defined, integrated units, reduced commercial fraud, and improved trade within the new nation. The conversion was defined in the United States as one meter equal to exactly 39.37 inches. This yields a conversion for a *U.S. survey foot* of:

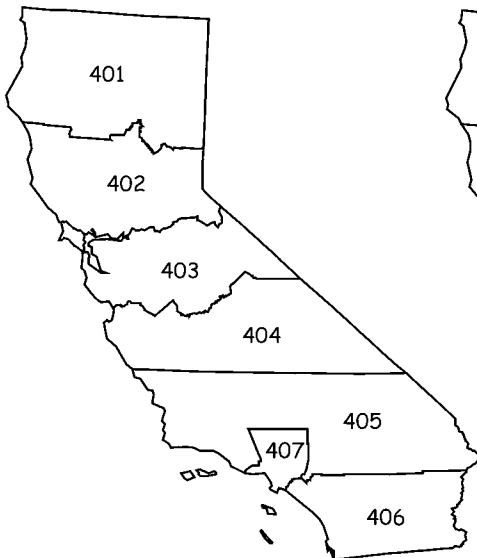
$$1 \text{ foot} = 0.3048006096012 \text{ meters}$$

Unfortunately, revolutionary tumult, national competition, and scientific differences led to the eventual adoption of a different conversion factor in Europe and most of the rest of the world. They adopted an *international foot* of:

$$1 \text{ foot} = 0.3048 \text{ meters}$$

The United States definition of a foot is slightly longer than the European definition,

FIPS codes of State Plane zones for use with NAD27



FIPS codes of State Plane zones for use with NAD83

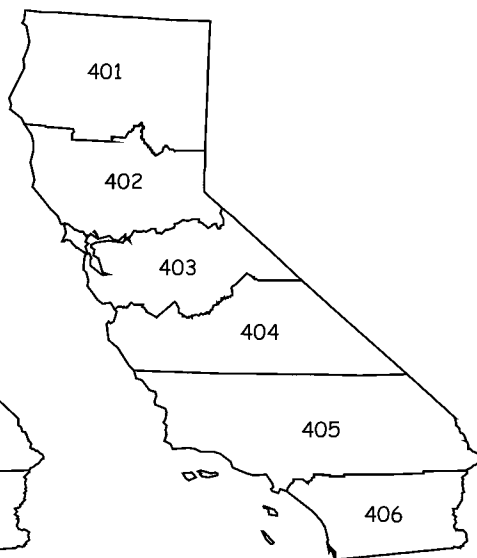


Figure 3-41: State Plane coordinate system zones and FIPS codes for California based on the NAD27 and NAD83 datums. Note that zone 407 from NAD27 is incorporated into zone 405 in NAD83.

by about one part in five million. Both conversions are used in the U.S., and the international conversion elsewhere. The European conversion was adopted as the standard for all measures under an international agreement in the 1950s. However, there was a long history of the use of the U.S. conversion in U.S. geodetic and land surveys. Therefore, the U.S. conversion was called the U.S. survey foot. This slightly longer metric-to-foot conversion factor should be used for all conversions among geodetic coordinate systems within the United States, for example, when converting from a State Plane coordinate system specified in feet to one specified in meters.

Universal Transverse Mercator Coordinate System

The Universal Transverse Mercator (UTM) coordinate system is another standard coordinate, distinct from the State Plane system. The UTM is a global coordinate system, based on the transverse Mercator projection. It is widely used in the United States and other parts of North America, and is also used in many other countries.

The UTM system divides the Earth into zones that are 6 degrees wide in longitude and extend from 80 degrees south latitude to 84 degrees north latitude. UTM zones are numbered from 1 to 60 in an easterly direction, starting at longitude 180 degrees West (Figure 3-42). Zones are further split north and south of the equator. Therefore, the zone containing most of England is identified as UTM Zone 30 North, while the zones containing most of New Zealand are designated UTM Zones 59 South and 60 South. Directional designations are here abbreviated, for example, 30N in place of 30 North.

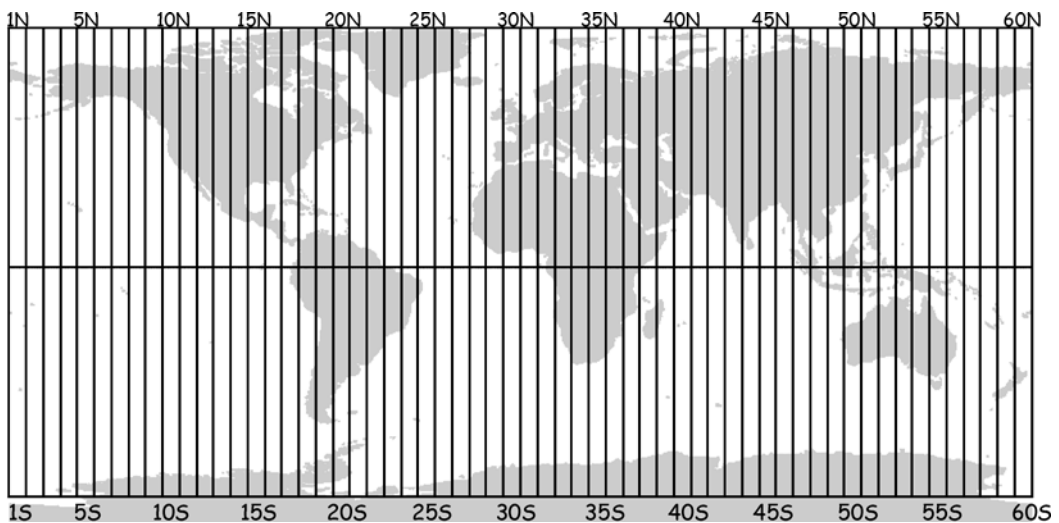


Figure 3-42: UTM zone boundaries and zone designators. Zones are six degrees wide and numbered from 1 to 60 from the International Date Line, 180°W. Zones are also identified by their position north and south of the equator, e.g., Zone 7 North, Zone 16 South.

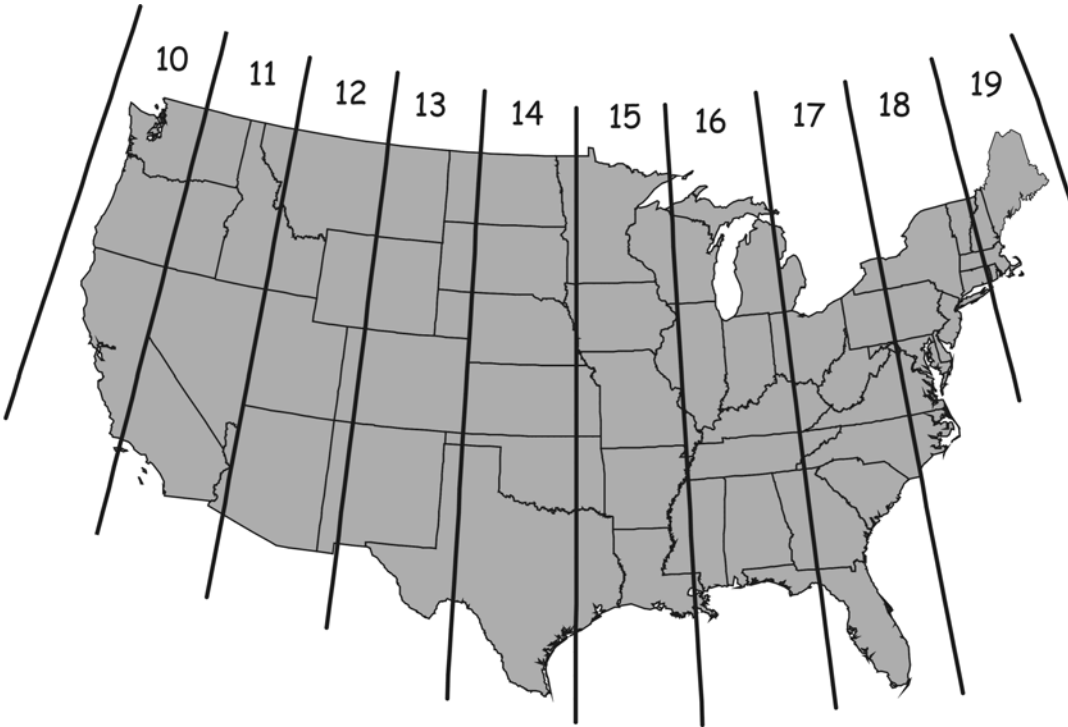


Figure 3-43: UTM zones for the lower 48 contiguous states of the United States of America. Each UTM zone is 6 degrees wide. All zones in the Northern Hemisphere are north zones, e.g., Zone 10 North, 11 North,...19 North.

The UTM coordinate system is common for data and study areas spanning large regions, for example, several State Plane zones. Many data from U.S. federal government sources are in a UTM coordinate system because many agencies manage large areas. Many state government agencies in the United States distribute data in UTM coordinate systems because the entire state fits predominantly or entirely into one UTM zone (Figure 3-43).

As noted before, all data for an analysis area must be in the same coordinate system if they are to be analyzed together. If not, the data will not co-occur as they should. The large width of the UTM zones accommodates many large-area analyses, and many states, national forests, or multi-county agencies have adopted the dominant UTM coordinate system as a standard.

We must note that the UTM coordinate system is not always compatible with

regional analyses. Because coordinate values are discontinuous across UTM zone boundaries, analyses are difficult across these boundaries. UTM zone 15 is a different coordinate system than UTM zone 16. The state of Wisconsin approximately straddles these two zones, and the state of Georgia straddles zones 16 and 17. If a uniform, statewide coordinate system is required, the choice of zone is not clear, and either one zone must be chosen, or some compromise projection must be chosen. For example, statewide analyses in Georgia and in Wisconsin are often conducted using UTM-like systems that involve moving the central meridian to near the center of each state.

Distances in the UTM system are specified in meters north and east of a zone origin (Figure 3-44). The y values are known as *northings*, and increase in a northerly direction. The x values are referred to as *eastings* and increase in an easterly direction.

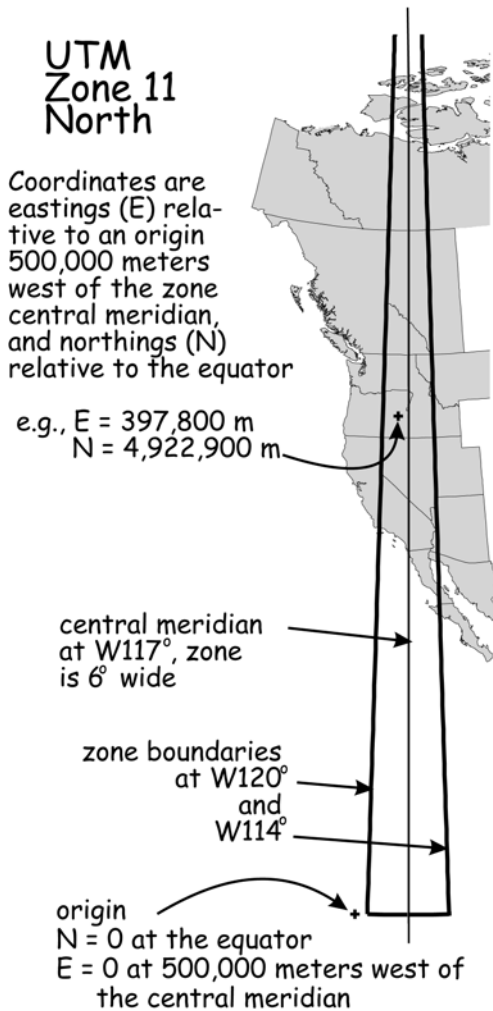


Figure 3-44: UTM zone 11N.

The origins of the UTM coordinate system are defined differently depending on whether the zone is north or south of the equator. In either case, the UTM coordinate system is defined so that all coordinates are positive within the zone. Zone easting coordinates are all greater than zero because the central meridian for each zone is assigned an easting value of 500,000 meters. This effectively places the origin ($E = 0$) at a point 500,000 meters west of the central meridian. All zones are less than 1,000,000 meters wide, ensuring that all eastings will be positive.

The equator is used as the northing origin for all north zones. Thus, the equator is assigned a northing value of zero for north zones. This avoids negative coordinates, because all of the UTM north zones are defined to be north of the equator.

University Transverse Mercator zones south of the equator are slightly different than those north of the equator (Figure 3-45). South zones have a *false northing* value added to ensure all coordinates within a zone are positive. UTM coordinate values increase as one moves from south to north in a projection area. If the origin were placed at the equator with a value of zero for south zone coordinate systems, then all the northing values would be negative. An offset is applied by assigning a false northing, a non-zero value, to an origin or other appropriate location. For UTM south zones, the northing values at the equator are set to equal 10,000,000 meters, assuring that all northing coordinate values will be positive within each UTM south zone (Figure 3-45).

Continental and Global Projections

There are map projections that are commonly used when depicting maps of continents, hemispheres, or other large regions. Just as with smaller areas, map projections for continental areas may be selected based on the distortion properties of the resultant map. Sizeable projection distortion in area, distance, and angle are observed in most large-area projections. Angles, distances, and areas are typically not measured or computed from these projections, as the differences between the map-derived and surface-measured values are too great for most uses. Large-area maps are most often used to display or communicate data for continental or global areas.

There are a number of projections that have been widely used for the world. These include variants of the Mercator, Goode, Mollweide, and Miller projections, among others. There is a trade-off that must be made in global projections, between a con-

UTM Zone 52 South

N=10,000,000 at the equator
E=0 at 500,000 meters west of
the central meridian

Coordinates are eastings (E) relative to an origin 500,000 meters west of the central meridian, and northings (N) relative to an origin 10,000,000 meters south of the equator

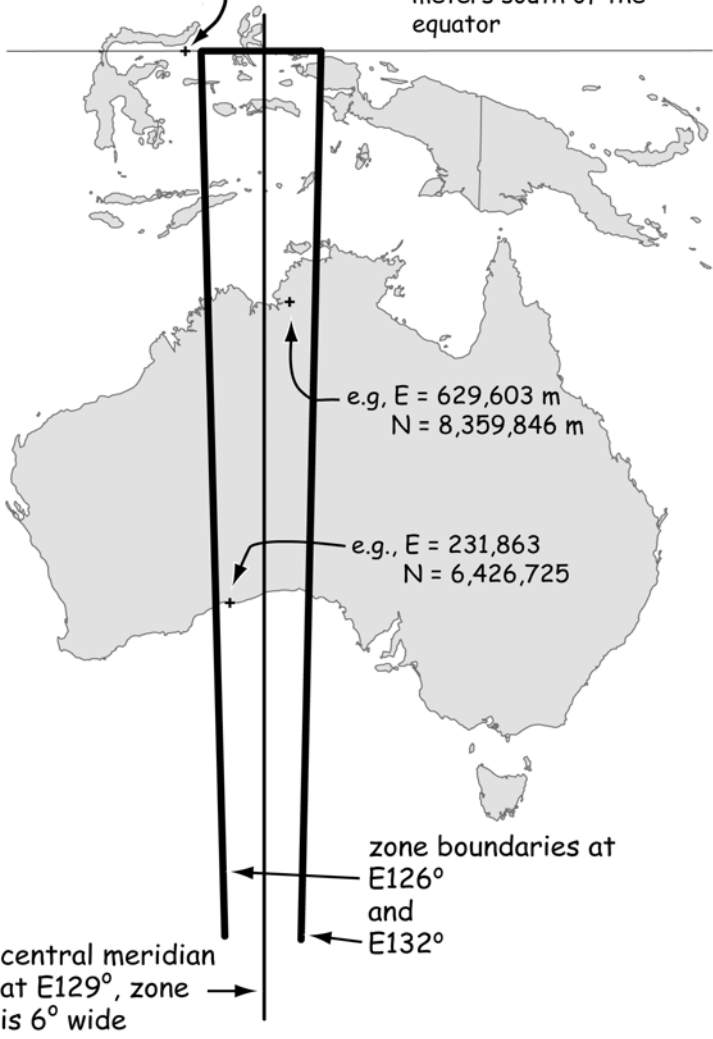


Figure 3-45: UTM south zones, such as Zone 52S shown here, are defined such that all the northing and easting values within the zone are positive. A false northing of 10,000,000 is applied to the equator, and a false easting of 500,000 is applied to the central meridian to ensure positive coordinate values throughout each zone.

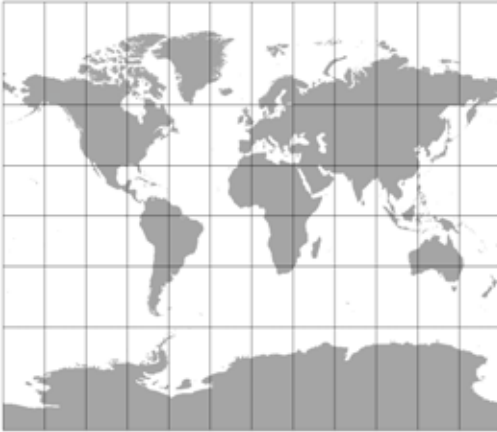


Figure 3-46: A Miller cylindrical projection, commonly used for maps of the world. This is an example of an uninterrupted map surface.

tinuous map surface and distortion. If a single, uncut surface is mapped, then there is severe distortion in some portion of the map. Figure 3-46 shows a Miller cylindrical projection, often used in maps of the world. This projection is similar to a Mercator projection, and is based on a cylinder that intersects the Earth at the equator. Distortion increases towards the poles, although not as much as with the Mercator.

Distortion in world maps may be reduced by using a cut or interrupted surface. Different projection parameters or sur-

faces may be specified for different parts of the globe. Projections may be mathematically constrained to be continuous across the area mapped.

Figure 3-47 illustrates an interrupted projection in the form of a Goode homolosine. This projection is based on a sinusoidal projection and a Mollweide projection. These two projection types are merged at parallels of identical scale. The parallel of identical scale in this example is set near the mid-northern latitude of $44^{\circ} 40' N$.

Continental projections may also be established. Generally, the projections are chosen to minimize area or shape distortion for the region to be mapped. Lambert conformal conic or other conic projections are often chosen for areas with a long east-west dimension, for example when mapping the contiguous 48 United States of America, or North America. Standard parallels are placed near the top and bottom of the continental area to reduce distortion across the region mapped. Transverse cylindrical projections are often used for large north-south continents.

None of these worldwide or continental projections are commonly used in a GIS for data storage or analysis. Uninterrupted coordinate systems show too much distortion to be of use in measuring most spatial quantities, and interrupted projections do not spec-

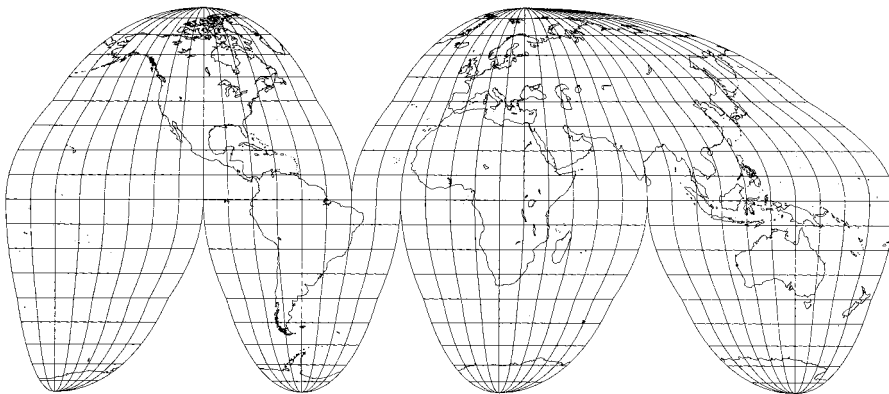


Figure 3-47: A Goode homolosine projection. This is an example of an interrupted projection, often used to reduce some forms of distortion when displaying the entire Earth surface. (from Snyder and Voxland, 1989)

ify a Cartesian coordinate system that defines positions for all points on the Earth surface. Worldwide data are typically stored in geographic coordinates (latitudes and longitudes). These data may then be projected to a specific coordinate system for display or document preparation.

Conversion Among Coordinate Systems

You might ask, how do I convert between geographic and projected coordinate systems? Exact or approximate mathematical formulas have been developed to convert to and from geographic latitude and longitude to all commonly used coordinate projections (Figure 3-48). These formulas are incorporated into “coordinate calculator” software packages, and are also integrated into most GIS software. For example, given a coordinate pair in the State Plane system, you may calculate the corresponding geographic coordinates. You may then apply a formula that converts geographic coordi-

nates to UTM coordinates for a specific zone using another set of equations. Since the backward and forward projections from geographic to projected coordinate systems are known, we may convert among most coordinate systems by passing through a geographic system (Figure 3-49, a).

Care must be taken when converting among projections that use different datums. If appropriate, we must insert a datum transformation when converting from one projected coordinate system to another (Figure 3-49, b). A datum transformation, described earlier in this chapter, is a calculation of the change in geographic coordinates when moving from one datum to another.

Users of GIS software should be careful when applying coordinate projection tools because the datum transformation may be omitted, or an inappropriate datum manually or automatically selected. For some software, the projection tool does not check or maintain information on the datum of the input spatial layer. This will often lead to an inappropriate or no datum transformation, and the output from the projection will be in error. Often these errors are small relative to other errors, for example, spatial imprecision in the collection of the line or point features. As shown in Figure 3-21, errors between NAD83(1986) and NAD83(CORS96) may be less than 10 cm (4 inches) in some regions, often much less than the average spatial error of the data themselves. However, errors due to ignoring the datum transformation may be quite large, for example, 10s to 100s of meters between NAD27 and most versions of NAD83, and errors of up to a meter are common between recent versions of WGS84 and NAD83. Given the sub-meter accuracy of many new GPS and other GNSS receivers used in data collection, datum transformation error of one meter is significant. As data collection accuracy improves, users develop applications based on those accuracies, so datum transformation errors should be avoided in all cases.

Conversion from geographic (lon, lat) to projected coordinates

Given longitude = λ , latitude = ϕ

Mercator projection coordinates are:

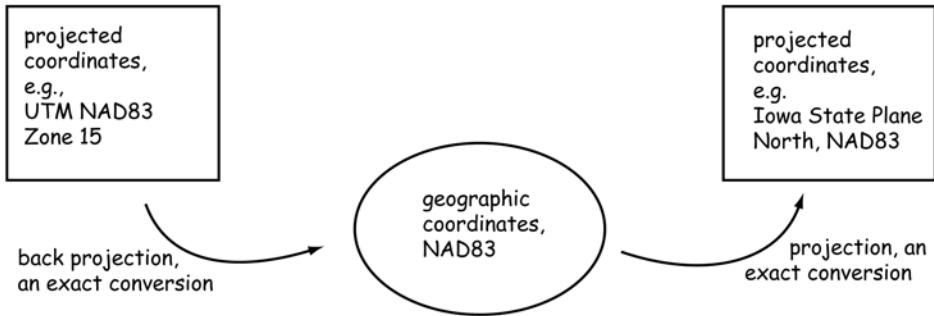
$$x = R \cdot (\lambda - \lambda_0)$$

$$y = R \cdot \ln (\tan (90^\circ + \phi/2))$$

where R is the radius of the sphere at map scale (e.g., Earth's radius), ln is the natural log function, and λ_0 is the longitudinal origin (Greenwich meridian)

Figure 3-48: Formulas are known for most projections that provide exact projected coordinates, if the latitudes and longitudes are known. This example shows the formulas defining the Mercator projection.

a) From one projection to another - same datum



b) From one projection to another - different datums

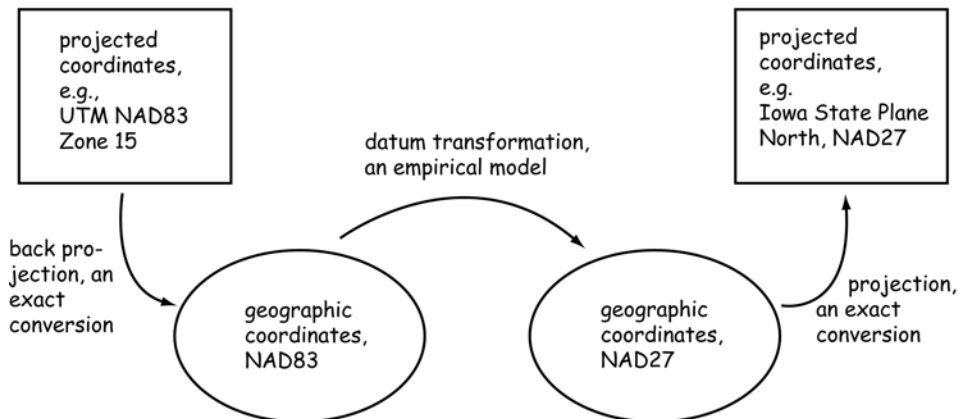


Figure 3-49: We may project between most coordinate systems via the back (or inverse) and forward projection equations. These calculate exact geographic coordinates from projected coordinates (a), and then new projected coordinates from the geographic coordinates. We must insert an extra step when a projection conversion includes a datum change. A datum transformation must be used to convert from one geodetic datum to another (b).

The Public Land Survey System

For the benefit of GIS practitioners in the United States we must cover one final land designation system, known as the *Public Land Survey System*, or PLSS. The PLSS is not a coordinate system, but PLSS points are often used as reference points in the United States, so the PLSS should be well understood for work there. The PLSS is a standardized method for designating and describing the location of land parcels. It was used for the initial surveys over most of the United States after the early 1800s, therefore nearly all land outside the original thirteen colonies uses the PLSS. An approximately uniform grid system was established across the landscape, with periodic adjustments incorporated to account for the anticipated error. Parcels were designated by their location within this grid system.

The PLSS was developed for a number of reasons. First, it was seen as a method to remedy many of the shortcomings of *metes and bounds* surveying, the most common method for surveying prior to the adoption of the PLSS. Metes and bounds describe a parcel relative to features on the landscape, sometimes supplemented with angle or distance measurements. In colonial times a parcel description might state “beginning at the joining of Shope Fork and Coweeta Creek, downstream along the creek approximately 280 feet to a large rock on the right bank, thence approximately northwest 420 feet to a large chestnut oak blazed with an S, thence west 800 feet to Mill Creek, thence down Mill Creek to Shope Fork Road, thence east on Shope Fork Road to the confluence of Shope Fork and Coweeta Creek.”

Metes and bounds descriptions require a minimum of surveying measurements, but it was a less than ideal system for describing locations or parcels. These metes and bounds descriptions could be vague, the features in the landscape might be moved or change, and it was difficult to describe parcels when there were few readily distinguished landscape features. Subdivided

6	5	4	3	2	1
7	8	9	10	11	12
18	17	16	15	14	13
19	20	21	22	23	24
30	29	28	27	26	25
31	32	33	34	35	36

Figure 3-50: Typical layout and section numbering of a PLSS township

parcels were often poorly described, and hence the source of much litigation, ill will, and many questionable real estate transactions.

The U.S. government needed a system that would provide unambiguous descriptions of parcels in unsettled territories west and south of the original colonies. The federal government saw public land sales as a way to generate revenue, to pay revolutionary war veterans, to expand the country, and to protect against encroachment by European powers. Parcels could not be sold until they were surveyed, therefore the PLSS was created. Land surveyed under the PLSS can be found in thirty states, including Alaska and most of the midwestern and western United States. Lands in the original 13 colonies, as well as West Virginia, Tennessee, Texas, and Kentucky were not surveyed under the PLSS system.

The PLSS divided lands by north-south lines running parallel to a principal meridian. New north-south lines were surveyed at six mile intervals. Additional lines were surveyed that were perpendicular to these north-south lines, in approximately east-west directions, and crossing meridian lines, also run at six-mile intervals. These lines form townships that were six miles square. Each township was further subdivided into 36 sections, each section approximately a

mile on a side. Each section was subdivided further, to quarter-sections (one-half mile on a side), or sixteenth sections, (one-quarter mile on a side, commonly referred to as quarter-quarter sections). Sections were numbered in a zig-zag pattern from one to 36, beginning in the northeast corner (Figure 3-50).

Surveyors typically marked the section corners and quarter-corners while running survey lines. Points were marked by a number of methods, including stone piles, pits, blaze marks chiseled in trees, and pipes or posts sunk in the ground.

Because the primary purpose of the PLSS survey was to identify parcels, lines and corner locations were considered static

on completion of the survey, even if the corners were far from their intended location. Survey errors were inevitable given the large areas and number of different survey parties involved. Rather than invite endless dispute and re-adjustment, the PLSS specifies that boundaries established by the appointed PLSS surveyors are unchangeable, and that township and section corners must be accepted as true. The typical section contains approximately 640 acres, but due in part to errors in surveying, sections larger than 1200 acres and smaller than 20 acres were also established (Figure 3-51).

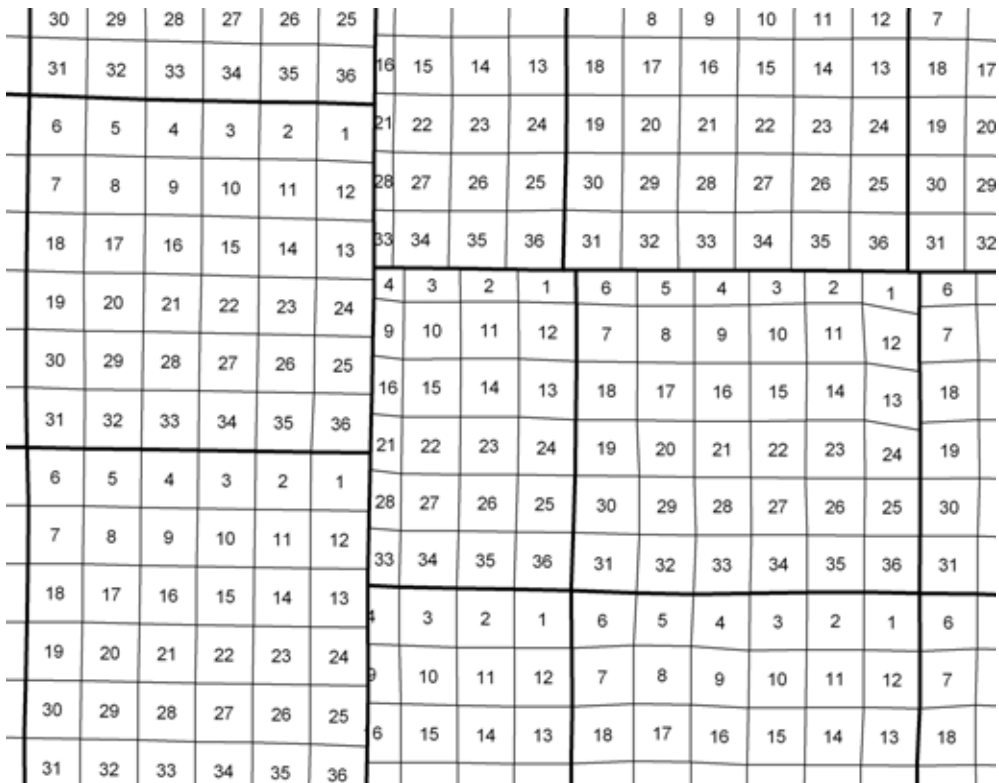


Figure 3-51: Example of variation in the size and shape of PLSS sections. Most sections are approximately one mile square with section lines parallel or perpendicular to the primary meridian, as illustrated by the township in the upper left of this figure. However, adjustments due to different primary meridians, different survey parties, and errors result in irregular section sizes and shapes.

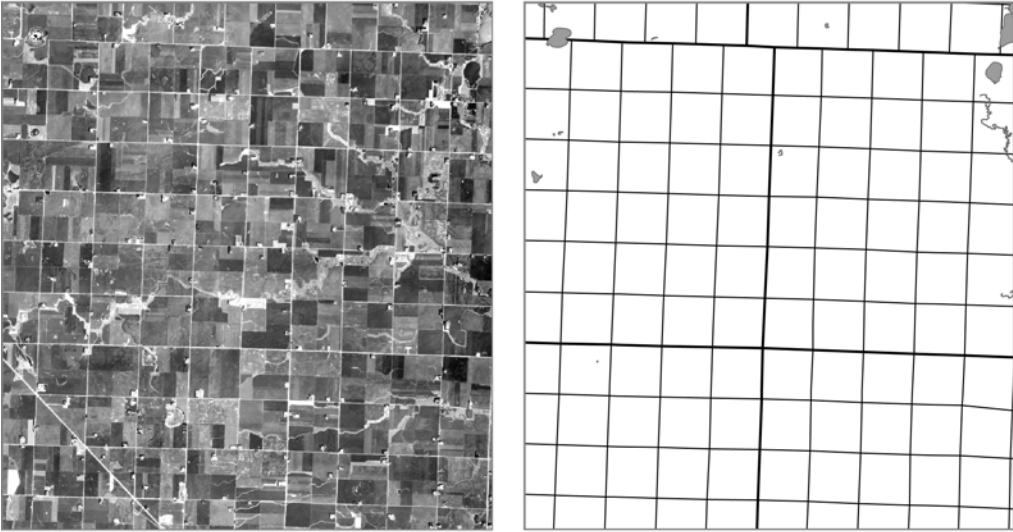


Figure 3-52: PLSS lines are often visible on the landscape. Roads (light lines on the image, above left) often follow the section and township lines (above right).

The PLSS is important today for several reasons. First, since PLSS lines are often property boundaries, they form natural corridors in which to place roads, powerlines, and other public services; they are often evident on the landscape (Figure 3-52). Many road intersections occur at PLSS corner points, and these can be viewed and referenced on many maps or imagery used for GIS database development efforts. Thus the PLSS often forms a convenient system to co-register GIS data layers. PLSS corners and lines are often plotted on government maps (e.g., 1:24,000 quads) or available as digital data (e.g., National Cartographic Information Center Digital Line Graphs). Further, PLSS corners are sometimes re-surveyed using high precision methods to provide property line control, particularly when a GIS is to be developed (Figure 3-53). Thus these points may be useful to properly locate and orient spatial data layers on the Earth's surface.



Figure 3-53: A PLSS corner that has been surveyed and marked with a monument. This monument shows the physical location of a section corner. These points are often used as control points for further spatial data development.

Summary

In order to enter coordinates in a GIS, we need to uniquely define the location of all points on Earth. We must develop a reference frame for our coordinate system, and locate positions on this system. Since the Earth is a curved surface and we work with flat maps, we must somehow reconcile these two views of the world. We define positions on the globe via geodesy and surveying. We convert these locations to flat surfaces via map projections.

We begin by modeling the Earth's shape with an ellipsoid. An ellipsoid differs from the geoid, a gravitationally-defined Earth surface, and these differences caused some early confusion in the adoption of standard global ellipsoids. There is a long history of ellipsoidal measurement, and we have arrived at our best estimates of global and regional ellipsoids after collecting large, painstakingly-developed sets of precise surface and astronomical measurements. These measurements are combined into datums, and these datums are used to specify the coordinate locations of points on the surface of the Earth.

Map projections are a systematic rendering of points from the curved Earth surface onto a flat map surface. While there are many purely mathematical or purely empirical map projections, the most common map projections used in GIS are based on developable surfaces. Cones, cylinders, and planes are the most common developable surfaces. A map projection is constructed by passing rays from a projection center through both the Earth surface and the developable surface. Points on the Earth are projected along the rays and onto the developable surface. This surface is then mathematically unrolled to form a flat map.

Standard sets of projections are commonly used for spatial data in a GIS. In the United States, the UTM and State Plane coordinate systems define a standard set of map projections that are widely used. Other map projections are commonly used for con-

tinental or global maps, and for smaller maps in other regions of the world.

A datum transformation is often required when performing map projections. Datum transformations account for differences in geographic coordinates due to changes in the shape or origin of the spheroid, and in some cases to datum adjustments. Datum transformation should be applied as a step in the map projection process when input and output datums differ.

A system of land division was established in the United States known as the Public Land Survey System (PLSS). This is not a coordinate system but rather a method for unambiguously and systematically defining parcels of land based on regularly spaced survey lines in approximately north-south and east-west directions. Intersection coordinates have been precisely measured for many of these survey lines, and are often used as a reference grid for further surveys or land subdivision.

Suggested Reading

- Bossler, J.D. (2002). Datums and geodetic systems, In J. Bossler (Ed.), *Manual of Geospatial Technology*. Taylor and Francis: London.
- Brandenburger, A.J. & Gosh, S K. (1985). The world's topographic and cadastral mapping operations. *Photogrammetric Engineering and Remote Sensing*, 51:437-444.
- Burkholder, E.F. (1993). Computation of horizontal/level distances. *Journal of Surveying Engineering*, 117:104-119.
- Colvocoresses, A.P. (1997). The gridded map. *Photogrammetric Engineering and Remote Sensing*, 63:371-376.
- Doyle, F.J. (1997). Map conversion and the UTM Grid. *Photogrammetric Engineering and Remote Sensing*, 63:367-370.
- Featherstone, W.E., & Kuhn, M. (2006). Height systems and vertical datums: a review in the Australian context. *Journal of Spatial Science*, 51:21-41.
- Habib, A. (2002). Coordinate transformation. In J. Bossler (Ed.), *Manual of Geospatial Technology*. Taylor and Francis: London.
- Flacke, W., & Kraus, B. (2005). *Working with Projections and Datum Transformations in ArcGIS: Theory and Practical Examples*. Points Verlag: Norden.
- Iliffe, J.C., & Lott, R. (2008). *Datums and Map Projections for Remote Sensing, GIS, and Surveying, 2nd ed.* CRC Press: Boca Raton.
- Janssen, V. (2009). Understanding coordinate reference systems, datums, and transformations. *International Journal of Geoinformatics*, 5:41-53.
- Keay, J. (2000). *The Great Arc*. Harper Collins: New York.
- Leick, A. (1993). Accuracy standards for modern three-dimensional geodetic networks. *Surveying and Land Information Systems*, 53:111-127.
- Maling, D.H. (1992). *Coordinate Systems and Map Projections*. George Phillip: London.
- Milbert, D. (2008). An analysis of the NAD83(NSRS2007) National Readjustment. Downloaded 9/12/2011 from http://www.ngs.noaa.gov/PUBS_LIB/NSRS2007

National Geospatial-Intelligence Agency (NGA), TR8350.2 World Geodetic System 1984, Its Definition and Relationship with Local Geodetic Systems. http://earth-info.nga.mil/GandG/publications/tr8350.2/tr8350_2.html

NOAA Manual NOS NGS 5 State Plane Coordinate System of 1983 -- http://www.ngs.noaa.gov/PUBS_LIB/ManualNOSNGS5.pdf

Schwartz, C.R. (1989). *North American Datum of 1983, NOAA Professional Paper NOS 2*. National Geodetic Survey: Rockville.

Smith, J. (1997). *Introduction to Geodesy: The History and Concepts of Modern Geodesy*, Wiley: New York.

Sobel, D. (1995). *Longitude*. Penguin Books: New York.

Soler, T. & Snay, R.A.(2004). Transforming positions and velocities between the International Terrestrial Reference Frame of 2000 and the North American Datum of 1983. *Journal of Surveying Engineering*, 130:49-55.

Snay, R.A. & Soler, T. (1999). Modern terrestrial reference systems, part 1. *Professional Surveyor*, 19:32-33.

Snay, R.A. & Soler, T. (2000). Modern terrestrial reference systems, part 2. the evolution of NAD83, *Professional Surveyor*, 20:16-18.

Snay, R.A., & Soler, T. (2000). Modern terrestrial reference systems, part 3. WGS84 and ITRS, *Professional Surveyor*, 20:24-28.

Snay, R.A., & Soler, T. (2000). Modern terrestrial reference systems, part 4, practical considerations for accurate positioning. *Professional Surveyor*, 20:32-34.

Snyder, J. (1993). *Flattening the Earth: Two Thousand Years of Map Projections*. Chicago: University of Chicago Press, Chicago.

Snyder, J. P. (1987). *Map Projections, A Working Manual, USGS Professional Paper No. 1396*. United States Government Printing Office: Washington D.C.

Snyder, J.P., & Voxland, P.M. (1989). *An Album of Map Projections, USGS Professional Paper No. 1453*. United States Government Printing Office: Washington D.C.

Tobler, W.R. (1962). A classification of map projections. *Annals of the Association of American Geographers*, 52:167-175.

U.S. Coast and Geodetic Survey Special Publication 235 The State Coordinate Systems -- http://www.ngs.noaa.gov/PUBS_LIB/publication235.pdf

- Van Sickle, J. (2010). *Basic GIS Coordinates, 2nd Edition*. CRC Press: Boca Raton.
- Vanicek, P., & Steevens, R.H. (1996). Transformation of coordinates between two horizontal geodetic datums. *Journal of Geodesy*, 70:740-745.
- Welch, R., & Homsey, A. (1997). Datum shifts for UTM coordinates. *Photogrammetric Engineering and Remote Sensing*, 63:371-376.
- Wolf, P. R., & Ghilani, C.D. (2002). *Elementary Surveying (10th ed.)*. Prentice-Hall: Upper Saddle River.
- Yang, Q., Snyder, J.P. & Tobler, W.R. (2000). *Map Projection Transformation: Principles and Applications*. Taylor & Francis: London.
- Zilkoski, D., Richards, J. & Young, G. (1992). Results of the general adjustment of the North American Vertical Datum of 1988. *Surveying and Land Information Systems*, 53:133-149.

Study Questions

3.1 - Can you describe how Eratosthenes estimated the circumference of the Earth? What value did he obtain?

3.2 - Assume the Earth is approximately a sphere (not an ellipsoid). Also assume you've repeated Poseidonius' measurements. What is your estimate of the radius of the Earth's sphere given the following distance/angle pairs. Note that the distances are given below meters, and angle in degrees, and calculators or spreadsheets may require you enter trigonometric angles in radians for trigonometric functions (1 radian = 57.2957795 degrees):

- a) angle = $1^{\circ} 18' 45.79558''$, distance = 146,000 meters
- b) angle = $0^{\circ} 43' 32.17917''$, distance = 80,500 meters
- c) angle = $0^{\circ} 3' 15.06032''$, distance = 6,000 meters

3.3 - Assume the Earth is approximately a sphere (not an ellipsoid). Also assume you've repeated Poseidonius' measurements. What is your estimate of the radius of the Earth's sphere given the following distance/angle pairs. Note that the distances are given below meters, and angle in degrees, and calculators or spreadsheets may require you enter trigonometric angles in radians for trigonometric functions (1 radian = 57.2957795 degrees):

- a) angle = $2^{\circ} 59' 31.33325''$, distance = 332,000 meters
- b) angle = $9^{\circ} 12' 12.77201''$, distance = 1,020,708 meters
- c) angle = $1^{\circ} 2' 12.15566''$, distance = 115,200 meters

3.4 - What is an ellipsoid? How does an ellipse differ from a sphere? What is the equation for the flattening factor?

3.5 - Why do different ellipsoids have different radii? Can you provide three reasons?

3.6 - Can you define the geoid? How does it differ from the ellipsoid, or the surface of the Earth? How do we measure the position of the geoid?

3.7 - Can you define a parallel or meridian in a geographic coordinate system? Where do the "horizontal" and "vertical" zero lines occur?

3.8 - How does magnetic north differ from the geographic North Pole?

3.9 - Can you define a datum? Can you describe how datums are developed?

3.10 - Why are there multiple datums, even for the same place on Earth? Can you define what we mean when we say there is a datum shift?

3.11 - What is a triangulation survey, and a bench mark?

3.12 - Why do we not measure vertical heights relative to mean sea level anymore?

3.13 - What is the difference between an orthometric height and a dynamic height.

3.14 - Use the NADCON software available from the U.S. NOAA/NGS website (http://www.ngs.noaa.gov/TOOLS/program_descriptions.html) to fill the following table. Note that all of these points are in CONUS, and longitudes are west, but entered as positive numbers.

Pnt	NAD27		NAD83(86)		HPGN	
	latitude	longitude	latitude	longitude	latitude	longitude
1	32°44'15"	117°09'42"	32°44'15.1827"	117°09'45.1202"	32°44'15.1820"	117°09'45.1200"
2	47°27'55"	122°18'06"	47°27'54.3574"	122°18'10.4453"	47°27'54.3642"	122°18'10.4366"
3	43°07'59"	89°20'11"	43°07'58.9806"	89°20'11.4226"		
4	29°58'07"	95°21'31"	29°58'07.7975"	95°21'31.7705"		
5	40°00'00"	105°16'01"			39°59'59.9552"	105°16'02.9712"
6	24°33'30"	81°45'19"			24°33'31.5216"	81°45'18.3362"
7			38°51'10.4052"	77°02'19.9165"	38°51'10.4063"	77°02'19.9041"
8			46°52'0.1524"	68°00'59.0974"	46°52'0.1580"	68°00'59.0995"

3.15 - Use the web version or download and start the HTDP software from the U.S. NOAA/NGS site listed above, and complete the following table. Enter epoch start and stop dates of 1, 1, 1986 and 1, 1, 2005, respectively, and specify a zero height or Z.

	NAD83(CORS96)		WGS(G1150)		ITRF2005	
Pnt	latitude	longitude	latitude	longitude	latitude	longitude
1	32°44'15"	117°09'42"	32°44'15.0321"	117°09'42.0662"	32°44'15.0325"	117°09'45.0663"
2	47°27'55"	122°18'06"	47°27'55.0183"	122°18'06.0583"	47°27'55.0186"	122°18'06.0583"
3	43°07'59"	89°20'11"	43°07'59.0283"	89°20'11.0293"		
4	29°58'07"	95°21'31"	29°58'07.0177"	95°21'31.0293"		
5	40°00'00"	105°16'01"			40°00'00.2143"	105°16'01.422"
6	24°33'30"	81°45'19"			24°33'30.0164"	81°45'19.0145"
7			38°51'1.0288"	77°02'21.0137"	38°51'1.0293"	77°02'21.0136"
8			46°52'0.0363"	68°01'00.0061"	46°52'0.0367"	68°01'01.0061"

3.16 - What is a developable surface? What are the most common shapes for a developable surface?

3.17 - Look up the NGS control sheets for the following points, and record their datums, latitudes and longitudes:

DOG, Maine, PID= PD0617.

Key West GSL, Florida, PID=AA1645

Neah A, Washington, PID=AF882

3.18 - Calculate the great circle distance for the control points, above, from:

- DOG to Neah A
- Key West to DOG
- Neah A to Key West

3.19 - Can you describe the State Plane coordinate system? What type of projections are used in a State Plane coordinate system?

3.20 - Can you define and describe the Universal Transverse Mercator coordinate system? What type of developable surface is used with a UTM projection? What are UTM zones, where is the origin of a zone, and how are negative coordinates avoided?

3.21 - What is a datum transformation? How does it differ from a map projection?

3.22 - Specify which type of map projection you would choose for each country, assuming you could use only one map projection for the entire country, the projection lines of intersection would be optimally-placed, and you wanted to minimize overall spatial distance distortion for the country. Choose from a transverse Mercator, a Lambert conformal conic, or an Azimuthal:

Benin	Bhutan
Slovenia	Israel

3.23 - Specify which type of map projection you would choose for each country, assuming you could use only one map projection for the entire country, the projection lines of intersection would be optimally-placed, and you wanted to minimize overall spatial distance distortion for the country. Choose from a transverse Mercator, a Lambert conformal conic, or an Azimuthal:

Chile	Nepal
Kyrgyzstan	The Gambia

3.24 - Can you describe the Public Land Survey System? Is it a coordinate system? What is its main purpose?