



Sampling The World

presented by:

Tim Haithcoat

**University of Missouri
Columbia**



Compiled with materials from:

Charles Parson, Bemidji State University
and

Timothy Nyerges, University of Washington

Introduction

- The world is infinitely complex
- The contents of a spatial database represent a particular view of the world
- The user sees the real world through the medium of the database
 - » The measurements & samples contained in the database must present as complete and accurate a view of the world as possible
 - » The contents of the database must be relevant in terms of:
 - themes and characteristics
 - The time period covered
 - The study area
- This presentation looks at techniques for sampling the world, and associated issues of accuracy, standards

Representing Reality

- A database consists of digital representations of discrete objects
- The features shown on a map (lakes, benchmarks, contours) can be thought of as discrete objects
 - » Thus the contents of a map can be captured in a database by turning map features into database objects
- Many of the features shown on a map are fictitious and do NOT exist in the real world
 - » Contours do not really exist, but houses & lakes are real objects
- The contents of a spatial database include:
 - » Digital versions of real objects (e.g. houses)
 - » Digital versions of artificial map features (e.g. contours)
 - » Artificial objects created for the purposes of the database (e.g. pixels)

Continuous Variation

- Some characteristics exist everywhere and vary continuously over the earth's surface
 - » Examples: elevation, atmospheric temperature and pressure, natural vegetation or soil type
- We can represent such variation in several ways:
 - » By taking measurements at sample points (e.g. weather stations)
 - » By taking transects
 - » By dividing the area into patches or zones, and assuming the variable is constant within each zone (e.g. soil mapping)
 - » By drawing contours (e.g. topographic mapping)

Continuous Variation (Continued)

- Each of these methods creates discrete objects
 - » The objects in each case are points, lines or areas
- A raster can be thought of as:
 - » A special case of a point sample where the points are regularly spaced
 - » A special case of zones where the zones are all the same size

Continuous Variation (Continued)

- Each method is approximate, capturing only part of the real variation
 - » A point sample misses variation between points
 - » Transects miss variation not on transects
 - » Zones pretend that variation is sudden at boundaries, and that there is no variation within zones
 - » Contours miss variation not located on contours
- Several methods can be used to try to improve the success of each method
 - » Example: for zones, map the boundaries as fuzzy instead of sharp lines
 - » Describe the zones as mixtures instead of as single classes, (e.g. 70% soil type A, 30% soil type B)

SPATIAL DATA

- Phenomena in the real world can be observed in three “modes”:

Spatial Mode

Deals with variation from place to place

Temporal Mode

Deals with variation from time to time

Thematic Mode

Deals with variation from one characteristic to another

Spatial Data (continued)

- All measurable or describable properties of the world can be considered to fall into one of these modes - place, time, and theme
- An exhaustive description of all three modes is not possible
- When observing real-world phenomena we usually hold one mode “fixed”, vary one in a “controlled” manner, and “measure” the third
 - » Example: using a census of population we could fix a time such as 1990, control for location using census tracts and measure a theme such as the percentage of persons owning automobiles

Spatial Data (continued)

- Holding geography fixed and varying time gives **longitudinal** data
- Holding time fixed and varying geography gives **cross-sectional** data
- The modes of information stored in a database influence the types of problem solving that can be accomplished
- The spatial mode of information is generally called **location**

Spatial Data (continued)

- **Attributes** capture the thematic mode by defining different characteristics of objects
- A table showing the attributes of objects is called an **attribute table**
 - » Each object corresponds to a row of the table
 - » Each characteristic or theme corresponds to a column of the table
 - » Thus, the table shows the thematic and some of the spatial modes

Spatial Data (continued)

- The temporal mode can be captured in several ways
 - » By specifying the interval of time over which an object exists
 - » By capturing information at certain points in time
 - » By specifying the rates of movement of objects
- Depending on how the temporal mode is captured, it may be included in a single attribute table, or be represented by a series of attribute tables on the same objects through time

SCALES OF MEASUREMENT

- Numerical values may be defined with respect to nominal, ordinal, interval, or ratio scales of measurement
- It is important to recognize the scales of measurement used in GIS data as this determines the kinds of mathematical operations that can be performed on the data
- The different scales can be demonstrated using an example of a marathon race.



NOMINAL SCALE

- On a nominal scale, numbers merely establish identity
 - » Example: a phone number signifies only the unique identity of the phone
- In the race, the numbers issued to racers which are used to identify individuals are a nominal scale
 - » These identity numbers do not indicate any order or relative value in terms of the race outcome



ORDINAL SCALE

- On an ordinal scale, numbers establish order only
 - » Example: phone number 961-8224 is not more of anything than 961-8049, so phone numbers are NOT ordinal
- In the race, the finishing places of each racer, i.e., 1st place, 2nd place, 3rd place, are measured on an ordinal scale
 - » However, we do NOT know how much time difference there is between each racer



INTERVAL SCALE

- On an interval scale, the difference (interval) between numbers is meaningful, but the numbering scale does not start at zero
 - » Subtraction makes sense, but division does not
 - » Example: it makes sense to say that 20°C is 10 degrees warmer than 10°C, so Celsius temperature is an interval scale, but 20°C is not twice as warm as 10°C
 - » Example: it makes no sense to say that the phone number 968-0244 is 62195 more than 961-8049, so phone numbers are not measurements on an interval scale

INTERVAL SCALE (continued)

- In the race, the time of the day that each racer finished is measured on an interval scale
 - » If the racers finished at 9:10, 9:20, and 9:25, then racer one finished 10 minutes before racer two and the difference between racers 1 and 2 is twice that of the difference between racers 2 and 3
 - » However, the racer finishing at 9:10 did not finish twice as fast as the racer finishing at 18:20



RATIO SCALE

- On a ratio scale, measurement has an absolute zero and the difference between numbers is significant
 - » Division makes sense
 - » Example: it makes sense to say that a 50 kg person weights half as much as a 100 kg person, so weight in kg is on a ratio scale
 - » The zero point of weight is absolute but the zero point of the Celsius scale is not

RATIO SCALE (continued)

- In our race, the first place finisher finished in a time of 2:30, the second in 2:30, and the 450th place finisher took 5 hours
 - » The 450th finisher took twice as long as the first place finisher ($5/2.5=2$)
- Note these distinctions, though important, are not always clearly defined
 - » Is elevation interval or ration? If the local base level is 750 ft, is a mountain at 2000 feet twice as high as one at 1000 feet when viewed from the valley?



RATIO SCALE (continued)

- Many types of geographical data used in GIS applications are nominal or ordinal
 - » Values establish the order of classes, or their distinct identify, but rarely intervals or ratios
- Thus, you cannot:
 - » Multiply soil type 2 by soil type 3 and get soil type 6
 - » Divide urban area by the rank of a city to get a meaningful number
 - » Subtract suitability class 1 from suitability class 4 to get 3 of anything
- However, you can:
 - » Divide population by area (both ratio scales) and get population density
 - » Subtract elevation at point A from elevation at point B to get the difference of elevation

NOMINAL

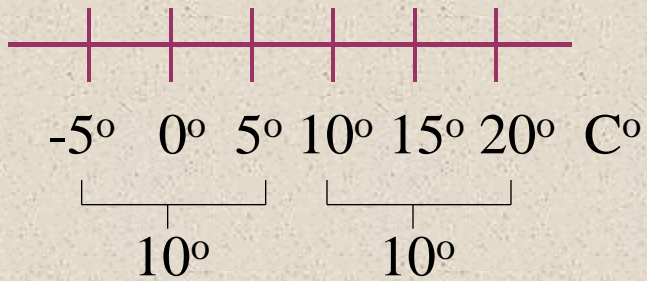
● Bill ● Mary
● Kim ● Rick

ORDINAL

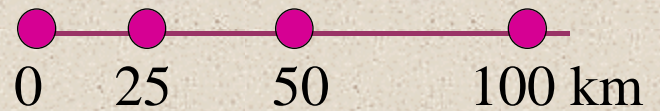
● ● ● ●
1st 2nd 3rd 4th

Scales of Measurement: Comparison

INTERVAL



RATIO

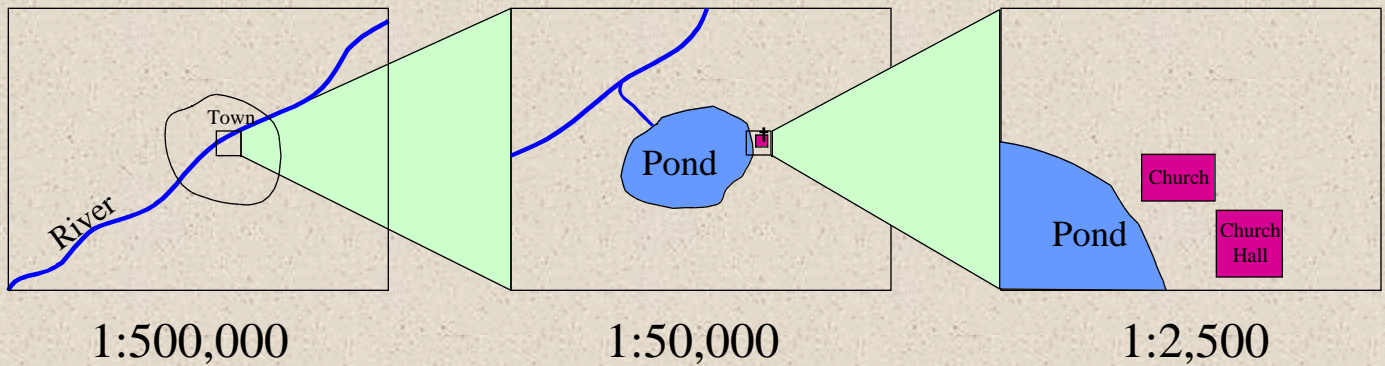


$$2 \times 50 \text{ km} = 100 \text{ km}$$

Multiple Representations

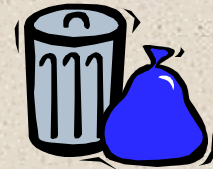
- A data model is essential to represent geographical data in a digital database
- There are many different data models
- The same phenomena may be represented in different ways, at different scales and with different levels of accuracy
- Thus, there may be multiple representations of the same geographical phenomena

Multiple Representations



Multiple Representations (continued)

- It is difficult to convert from one representation to another
 - » From a small scale (1:250,000) to a large scale (1:10,000)
- Thus, it is common to find databases with multiple representations of the same phenomenon
 - » This is wasteful, but techniques to avoid it are poorly developed



Primary Data Collection

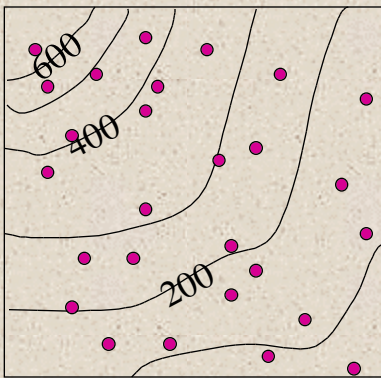
- Some data in a spatial database may have been measured directly
 - » Example: by field sampling or remote sensing
- The density of sampling determines the resolution of the data
 - » Example: samples taken every hour will capture hour-to-hour variation, but miss shorter-term variation
 - » Example: samples taken every 1 km will miss any variation at resolutions less than 1 km



Primary Data Collection (continued)

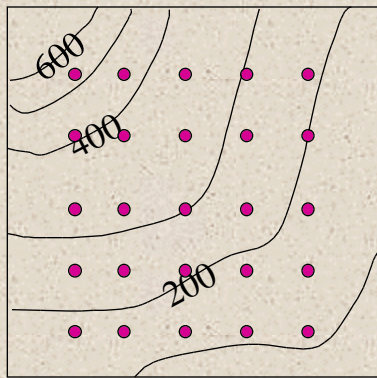
- A sample is designed to capture the variation present in a larger universe
 - » Example: a sample of places should capture the variation present at all possible places
 - » Example: a sample of times will be designed to capture variation at all possible times
- There are several approaches to sampling
 - » Random Samples
 - » Systematic Samples
 - » Stratified Samples

Sampling Strategies



Random

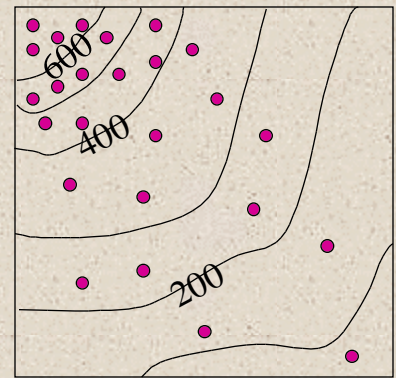
Every place or time is equally likely to be chosen.



Systematic

Samples are chosen according to a rule.

Example: every 1 km, but the rule is expected to create no bias in the results of analysis



Stratified

Research knows for some reason that the universe contains significantly different sub-populations, & samples within each sub-population in order to achieve adequate representation of each.

Secondary Data Sources

- Some data may have been obtained from existing maps, tables, or other databases
 - » Such sources are termed secondary
- To be useful, it is important to obtain information in addition to the data themselves:
 - » Information on the procedures used to collect and compile the data
 - » Information on coding schemes, accuracy of instruments
- Unfortunately, such information is often not available
 - » A user of a spatial database may not know how the data were captured and processed prior to input
 - » This often leads to misinterpretation, false expectations about accuracy

Standards

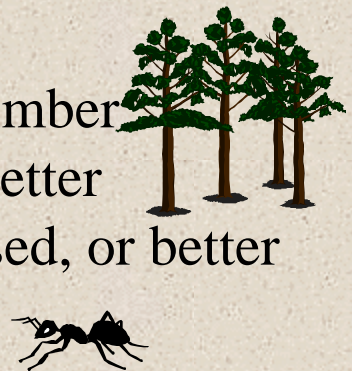
- Standards may be set to assure uniformity

- » Within a single data set
- » Across data sets

- » Example: uniform information about timber types throughout the database allows better



fire fighting methods to be used, or better control of insect infestations.





- Data capture should be undertaken in standardized ways that will assure the widest possible use of the information

Sharing Data

- It is not uncommon for as many as three agencies to create databases with, ostensibly, the same information
 - » Examples:
 - A planning agency may map land use, including a forested class
 - The state department of forestry also maps forests
 - The wildlife division of the department of conservation maps habitat, which includes fields and forest



Sharing Data (continued)

- Each may digitize their forest class onto different GIS systems, using different protocols, and with different definitions for the classes of forest cover
- This is a waste of time  and money 
- Sharing information gives it added value
- Sharing basic formats with other information providers, such as a department of transportation, might make marketing the database more profitable

Agency Standards

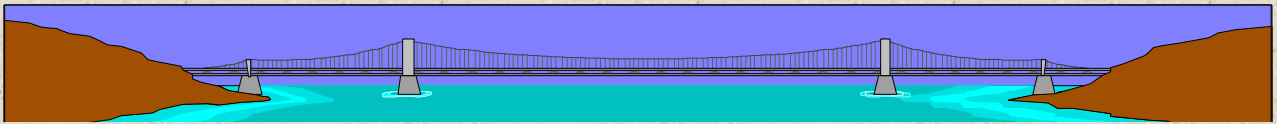
- State & national agencies have set standards for certain environmental data
 - » The Soil Conservation Service (SCS) has adopted the “seventh approximation” as the national taxonomy
 - » The US Geological Survey has set standards for landuse, transportation, and hydrography that are used as guidelines in many states
 - » Forest inventories are not standardized; agencies may use different systems while managing a contiguous region of forest land

Errors and Accuracy

- There's nearly a universal tendency to lose sight of errors once the data are in digital form
- Errors:
 - » Are implanted in original databases because of errors in the original sources (source errors)
 - » are added during data capture and storage (processing errors)
 - » Occur when data are extracted from the computer
 - » Arise when the various layers of data are combined in an analytical exercise

Original Sin: Source Errors

- Are extremely common in non-mapped source data, such as locations of wells, or lot descriptions
- Can be caused by doing inventory work from aerial photography and misinterpreting images
- Often occur because base maps are relied on too heavily



- A recent attempt in Minnesota to overlay Department of Transportation bridge locations on USGS transportation data resulted in bridges lying neither beneath roads, nor over water, and roads lying apparently under rivers
- Until they were compared in this way, it was assumed that each data set was locationally acceptable
- The ability of GIS to overlay may expose previously unsuspected errors

Boundaries



- Boundaries of soil types are actually transition zones, but are mapped by lines less than 0.5 mm wide

- Lakes fluctuate widely in area, yet have permanently recorded shorelines

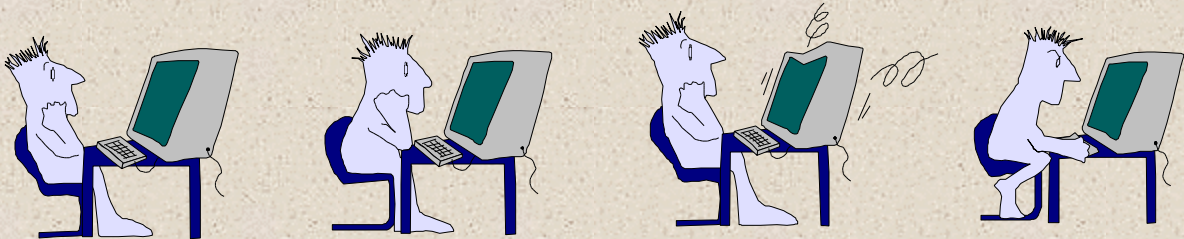


Classification Errors

- Common when tabular data are rendered in map form
- Simple typing errors may be invisible until presented graphically
 - » Floodplain soils may appear on hilltops
 - » Pastureland may appear to be misinterpreted marsh
- More complex classification errors may be due to sampling strategies that produce the original data
- Timber appraisal is commonly done using a few, randomly selected points to describe large stands
 - » Information may exist that documents the error of the sampling technique
 - » However, such information is seldom included in the GIS database

Data Capture Errors

- Manual data input induces another set of errors
- Eye-hand coordination varies from operator to operator and from time to time
 - » Data input is a tedious task, it is difficult to maintain quality over long periods of time



Accuracy Standards

- Many agencies have established accuracy standards for geographical data
 - » These are more often concerned with accuracy of locations of objects than with accuracy of attributes
- Location accuracy standards are commonly decided from the scale of source materials
 - » For natural resource data 1:24,000 scale accuracy is a common target
 - » At this scale, 0.5 mm line width = 12 m on the ground

Accuracy Standards (continued)

- USGS topographic information is currently available in digital form at 1:100,000
 - » 0.5 mm line width = 50 m on the ground
- Higher accuracy requires better source materials
 - » Is the added cost justified by the objectives of the study?
- Accuracy standards should be determined by considering both the value of information and the cost of collection

