

The third case above (spatial point pattern data) could be exemplified by residences of persons suffering from a particular disease, or by locations of a certain species of tree in a forest. Here the response  $Y$  is often fixed (occurrence of the event), and only the locations  $\mathbf{s}_i$  are thought of as random. In some cases this information might be supplemented by age or other covariate information, producing a *marked* point pattern). Such data are often of interest in studies of event *clustering*, where the goal is to determine whether an observed spatial point pattern is an example of a clustered process (where points tend to be spatially close to other points), or merely the result of a random event process operating independently and homogeneously over space. Note that in contrast to areal data, where no individual points in the data set could be identified, here (and in point-referenced data as well) precise locations are known, and so must often be protected to protect the privacy of the persons in the set.

Even though our preferred inferential outlook is Bayesian, the statistical inference tools discussed in Chapters 2 through 4 are entirely classical. While all subsequent chapters adopt the Bayesian point of view, our objective here is to acquaint the reader with the classical techniques first, since they are more often implemented in standard software packages. Moreover, as in other fields of data analysis, classical methods can be easier to compute, and produce perfectly acceptable results in relatively simple settings. Classical methods often have interpretations as limiting cases of Bayesian methods under increasingly vague prior assumptions. Finally, classical methods can provide insight for formulating and fitting hierarchical models.

In the case of point-level data, the location index  $\mathbf{s}$  varies *continuously* over  $D$ , a fixed subset of  $\mathbb{R}^d$ . Suppose we assume that the covariance between the random variables at two

locations depends on the *distance* between the locations. One frequently used association specification is the exponential model. Here the covariance between measurements at two locations is an exponential function of the interlocation distance, i.e.,  $Cov(Y(\mathbf{s}_i), Y(\mathbf{s}_{i'})) \equiv C(d_{ii'}) = \sigma^2 e^{-\phi d_{ii'}}$  for  $i \neq i'$ , where  $d_{ii'}$  is the distance between sites  $\mathbf{s}_i$  and  $\mathbf{s}_{i'}$ , and  $\sigma^2$  and  $\phi$  are positive parameters called the *partial sill* and the *decay parameter*, respectively ( $1/\phi$  is called the *range parameter*). A plot of the covariance versus distance is called the *covariogram*. When  $i = i'$ ,  $d_{ii'}$  is of course 0, and  $C(d_{ii'}) = Var(Y(\mathbf{s}_i))$  is often expanded to  $\tau^2 + \sigma^2$ , where  $\tau^2 > 0$  is called a *nugget effect*, and  $\tau^2 + \sigma^2$  is called the *sill*. Of course, while the exponential model is convenient and has some desirable properties, many other parametric models are commonly used; see Section 2.1 for further discussion of these and their relative merits.

Adding a joint distributional model to these variance and covariance assumptions then enables likelihood inference in the usual way. The most convenient approach would be to assume a multivariate *normal* (or *Gaussian*) distribution for the data. That is, suppose we are given observations  $\mathbf{Y} \equiv \{Y(\mathbf{s}_i)\}$  at known locations  $\mathbf{s}_i$ ,  $i = 1, \dots, n$ . We then assume that

$$\mathbf{Y} \mid \mu, \boldsymbol{\theta} \sim N_n(\mu \mathbf{1}, \Sigma(\boldsymbol{\theta})) \text{ ,} \quad (1.1)$$

where  $N_n$  denotes the  $n$ -dimensional normal distribution,  $\mu$  is the (constant) mean level,  $\mathbf{1}$  is a vector of ones, and  $(\Sigma(\boldsymbol{\theta}))_{ii'}$  gives the covariance between  $Y(\mathbf{s}_i)$  and  $Y(\mathbf{s}_{i'})$ . For the variance-covariance specification of the previous paragraph, we have  $\boldsymbol{\theta} = (\tau^2, \sigma^2, \phi)^T$ , since the covariance matrix depends on the nugget, sill, and range.

In fact, the simplest choices for  $\Sigma$  are those corresponding to *isotropic* covariance functions, where we assume that the spatial correlation is a function solely of the distance  $d_{ii'}$  between  $\mathbf{s}_i$  and  $\mathbf{s}_{i'}$ . As mentioned above, exponential forms are particularly intuitive examples. Here,

$$(\Sigma(\boldsymbol{\theta}))_{ii'} = \sigma^2 \exp(-\phi d_{ii'}) + \tau^2 I(i = i'), \quad \sigma^2 > 0, \phi > 0, \tau^2 > 0, \quad (1.2)$$

where  $I$  denotes the indicator function (i.e.,  $I(i = i') = 1$  if  $i = i'$ , and 0 otherwise). Many other choices are possible for  $Cov(Y(\mathbf{s}_i), Y(\mathbf{s}_{i'}))$ , including for example the powered exponential,

$$(\Sigma(\boldsymbol{\theta}))_{ii'} = \sigma^2 \exp(-\phi d_{ii'}^\kappa) + \tau^2 I(i = i'), \quad \sigma^2 > 0, \phi > 0, \tau^2 > 0, \kappa \in (0, 2],$$

the spherical, the Gaussian, and the Matérn (see Subsection 2.1.3 for a full discussion). In particular, while the latter requires calculation of a modified Bessel function, Stein (1999a, p. 51) illustrates its ability to capture a broader range of local correlation behavior despite having no more parameters than the powered exponential. We shall say much more about point-level spatial methods and models in Chapters 2, 3 and 6 and also provide illustrations using freely available statistical software.

### 1.1.2 Areal models

In models for areal data, the geographic regions or *blocks* (zip codes, counties, etc.) are denoted by  $B_i$ , and the data are typically sums or averages of variables over these blocks. To introduce spatial association, we define a *neighborhood* structure based on the arrangement of the blocks in the map. Once the neighborhood structure is defined, models resembling autoregressive time series models are considered. Two very popular models that incorporate such neighborhood information are the *simultaneously* and *conditionally autoregressive* models (abbreviated SAR and CAR), originally developed by Whittle (1954) and Besag (1974), respectively. The SAR model is computationally convenient for use with likelihood methods. By contrast, the CAR model is computationally convenient for Gibbs sampling used









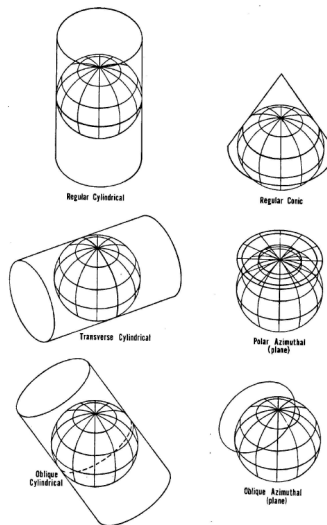


Figure 1.4 *The geometric constructions of projections using developable surfaces (figure courtesy of the U.S. Geological Survey).*

and the patches are closely approximated by planes) and deriving a set of (partial) differential equations whose solution will yield  $f$  and  $g$ . Suitable initial conditions are set to create projections with desired geometric properties.

Thus, consider a small patch on the sphere formed by the infinitesimal quadrilateral,  $ABCD$ , given by the vertices,

$$A = (\lambda, \phi), \quad B = (\lambda, \phi + d\phi), \quad C = (\lambda + d\lambda, \phi), \quad D = (\lambda + d\lambda, \phi + d\phi).$$

So, with  $R$  being the radius of the earth, the horizontal differential component along an arc of latitude is given by  $|AC| = (R \cos \phi)d\lambda$  and the vertical component along a great circle of longitude is given by  $|AB| = Rd\phi$ . Note that since  $AC$  and  $AB$  are arcs along the latitude and longitude of the globe, they intersect each other at right angles. Therefore, the area of the patch  $ABCD$  is given by  $|AC||AB|$ . Let  $A'B'C'D'$  be the (infinitesimal) image of the patch  $ABCD$  on the map. Then, we see that

$$\begin{aligned} A' &= (f(\lambda, \phi), g(\lambda, \phi)), \\ C' &= (f(\lambda + d\lambda, \phi), g(\lambda + d\lambda, \phi)), \\ B' &= (f(\lambda, \phi + d\phi), g(\lambda, \phi + d\phi)), \\ \text{and } D' &= (f(\lambda + d\lambda, \phi + d\phi), g(\lambda + d\lambda, \phi + d\phi)). \end{aligned}$$

This in turn implies that

$$\overrightarrow{A'C'} = \left( \frac{\partial f}{\partial \lambda}, \frac{\partial g}{\partial \lambda} \right) d\lambda \quad \text{and} \quad \overrightarrow{A'B'} = \left( \frac{\partial f}{\partial \phi}, \frac{\partial g}{\partial \phi} \right) d\phi.$$

If we desire an equal-area projection we need to equate the area of the patches  $ABCD$  and  $A'B'C'D'$ . But note that the area of  $A'B'C'D'$  is given by the area of parallelogram formed by vectors  $\overrightarrow{A'C'}$  and  $\overrightarrow{A'B'}$ . Treating them as vectors in the  $xy$  plane of an  $xyz$  system, we see that the area of  $A'B'C'D'$  is the cross-product,

$$(\overrightarrow{A'C'}, 0) \times (\overrightarrow{A'B'}, 0) = \left( \frac{\partial f}{\partial \lambda} \frac{\partial g}{\partial \phi} - \frac{\partial f}{\partial \phi} \frac{\partial g}{\partial \lambda} \right) d\lambda d\phi.$$



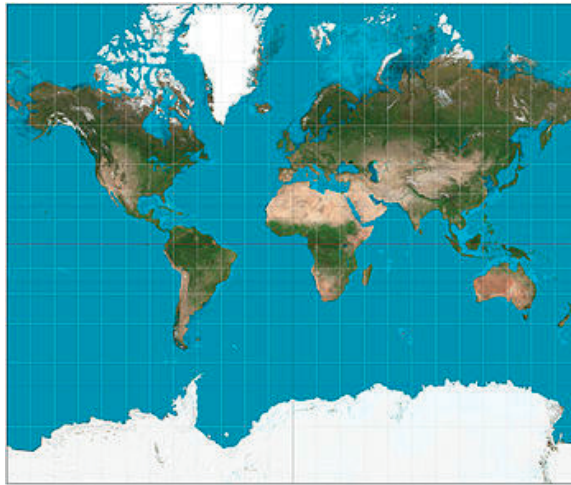


Figure 1.6 *The Mercator projection.*

$R \sec \phi$ . After suitable integration, this leads to the analytical equations (with the 0 degree meridian as the central meridian),

$$f(\lambda, \phi) = R\lambda; \quad g(\lambda, \phi) = R \ln \tan \left( \frac{\pi}{4} + \frac{\phi}{2} \right).$$

As is seen above, even the simplest map projections lead to complex transcendental equations relating latitude and longitude to positions of points on a given map. Therefore, rectangular grids have been developed for use by surveyors. In this way, each point may be designated merely by its distance from two perpendicular axes on a flat map. The  $y$ -axis usually coincides with a chosen central meridian,  $y$  increasing north, and the  $x$ -axis is perpendicular to the  $y$ -axis at a latitude of origin on the central meridian, with  $x$  increasing east. Frequently, the  $x$  and  $y$  coordinates are called “eastings” and “northings,” respectively, and to avoid negative coordinates, may have “false eastings” and “false northings” added to them. The grid lines usually do not coincide with any meridians and parallels except for the central meridian and the equator.

One such popular grid, adopted by The National Imagery and Mapping Agency (NIMA) (formerly known as the Defense Mapping Agency), and used especially for military use throughout the world, is the Universal Transverse Mercator (UTM) grid; see Figure 1.7. The UTM divides the world into 60 north-south zones, each of width six degrees longitude. Starting with Zone 1 (between 180 degrees and 174 degrees west longitude), these are numbered consecutively as they progress eastward to Zone 60, between 174 degrees and 180 degrees east longitude. Within each zone, coordinates are measured north and east in meters, with northing values being measured continuously from zero at the equator, in a northerly direction. Negative numbers for locations south of the equator are avoided by assigning an arbitrary false northing value of 10,000,000 meters (as done by NIMA's cartographers). A central meridian cutting through the center of each 6 degree zone is assigned an easting value of 500,000 meters, so that values to the west of the central meridian are less than 500,000 while those to the east are greater than 500,000. In particular, the conterminous 48 states of the United States are covered by 10 zones, from Zone 10 on the west coast through Zone 19 in New England.

In practice, the UTM is used by overlaying a transparent grid on the map, allowing distances to be measured in meters at the map scale between any map point and the nearest grid lines to the south and west. The northing of the point is calculated as the sum







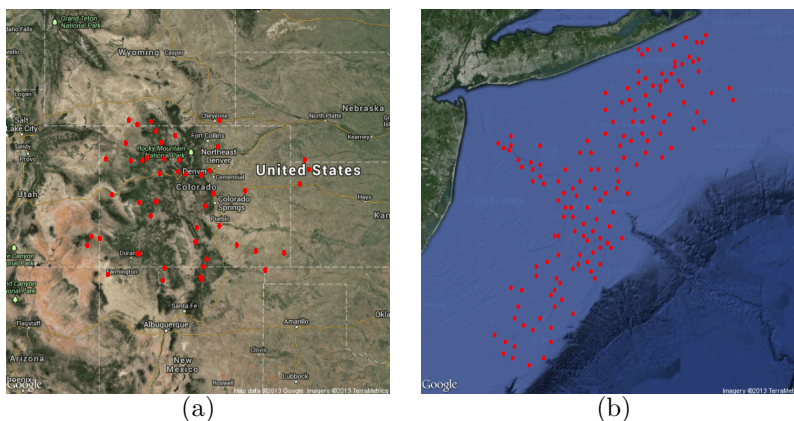












```
> SP_longlat <- SpatialPoints(coords =
+                               cbind(coloradoST$Longitude,
+                                     coloradoST$Latitude),
+                               proj4string = CRS("+proj=longlat +ellps=WGS84"))
> SP_utm <- spTransform(SP_longlat,
+                        CRS("+proj=utm +zone=13 +datum=WGS84"))
> plot(SP_utm).
```

Finally, we mention the `fields` package in R, which offers several useful functions for spatial analysis. In particular, it includes two functions `rdist` and `rdist.earth` that conveniently compute inter-site distances. Let `X1` and `X2` be two matrices representing two different sets of locations. Then,

computes the inter-site distance matrices between the locations in **X1** and **X2**. The function **rdist** uses the Euclidean distance, while **rdist.earth** uses the spherical or geodetic distance. The latter should be used only when **X1** and **X2** contain latitude-longitude coordinates.

1. What sorts of areal unit variables can you envision that could be viewed as arising from point-referenced variables? What sorts of areal unit variables can you envision whose mean could be viewed as arising from a point-referenced surface? What sorts of areal unit variables fit neither of these scenarios?



(a) Compute the above projection for Chicago and Minneapolis ( $N = 2$ ) and find the Euclidean distance between the projected coordinates. Compare with the geodesic distance. Repeat this exercise for New York and New Orleans.

11. Use the `sp`, `rgdal` and `RgoogleMaps` packages to create an UTM projection for the locations in the scallops data and produce the picture in Figure 1.11(b).

12. Use the **fields** package to produce the inter-site distance matrix for the locations in the scallops data. Compute this matrix using the **rdist.earth** function, which yields the geodetic distances. Next project the data to UTM coordinates and use the **rdist** function to compute the inter-site Euclidean distance matrix. Draw histograms of the inter-site distances and comment on any notable discrepancies resulting from the map projection.