

1

Introduction

1.1 Background

Ocean modelers, in the formal sense of the word, attempting to describe the ocean circulation have paid comparatively little attention in the past to the problems of working with real data. Thus, for example, one drives models, theoretical or numerical, with analytically prescribed wind or buoyancy forcing without worrying overly much about how realistic such assumed forms might be. The reasons for approaching the problem this way are good ones—there has been much to learn about how the models themselves behave, without troubling initially about the question of whether they describe the real ocean. Furthermore, there has been extremely little in the way of data available, even had one wished to use, say, realistic wind and buoyancy flux fields.

This situation is changing rapidly; the advent of wind measurements from satellite-borne instruments and other improvements in the ability of meteorologists to estimate the windfields over the open ocean, and the development of novel technologies for observing the ocean circulation, have made it possible to seriously consider estimating the global circulation in ways that were visionary only a decade ago. Technologies of neutrally buoyant floats, long-lived current meters, chemical tracer observations, satellite altimeters, acoustical methods, etc., are all either here or imminent.

The models themselves have also become so complex (e.g., Figure 1–1a) that special tools are required to understand them, to determine whether they are actually more complex than required to describe what we see (Figure 1–1b), or if less so, to what externally prescribed parameters or missing internal physics they are likely to be most sensitive.

The ocean is so difficult to observe that theoreticians intent upon explaining known phenomena have made plausible assumptions about the behavior

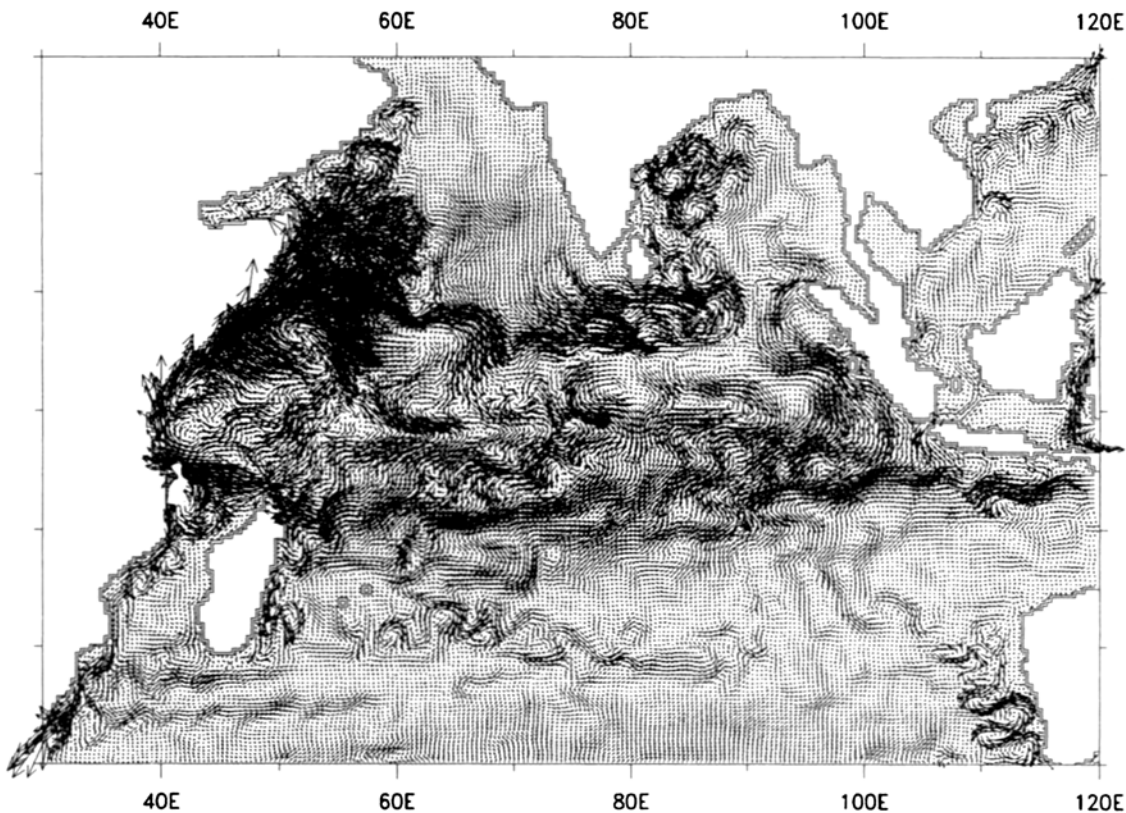


Figure 1–1a. The sear-surface velocity on 15 September 1988 as estimated from a high-resolution global, oceanic general circulation model. Only a fraction of the global domain is depicted.

The model is a nominal $1/4^\circ$ latitude-longitude resolution version of the earlier computation described by Semtner and Chervin (1992). Much of the visual structure is time-dependent and raises very

serious issues of observational sampling and of mathematical representation. (Courtesy of R. Tokmakian and A. Semtner.)

of the system and proceeded to construct systems of equations that are then solved. Some of these assumptions are indeed so plausible, and the resulting calculations so interesting, that it has become forgotten that they were assumed rather than demonstrated. The assumptions become elevated to the level of textbook dogma as known. One thinks immediately of the use of eddy coefficients in Laplacian diffusion/mixing terms, the assumption that Ekman layer divergences drive the large-scale interior circulation, that this circulation is in Sverdrup balance, that steady models are adequate descriptors of the circulation, etc. Consequently, the subject is rife with myths. Myths are important elements in human views of the world and often contain major components of historical truth. They have much that

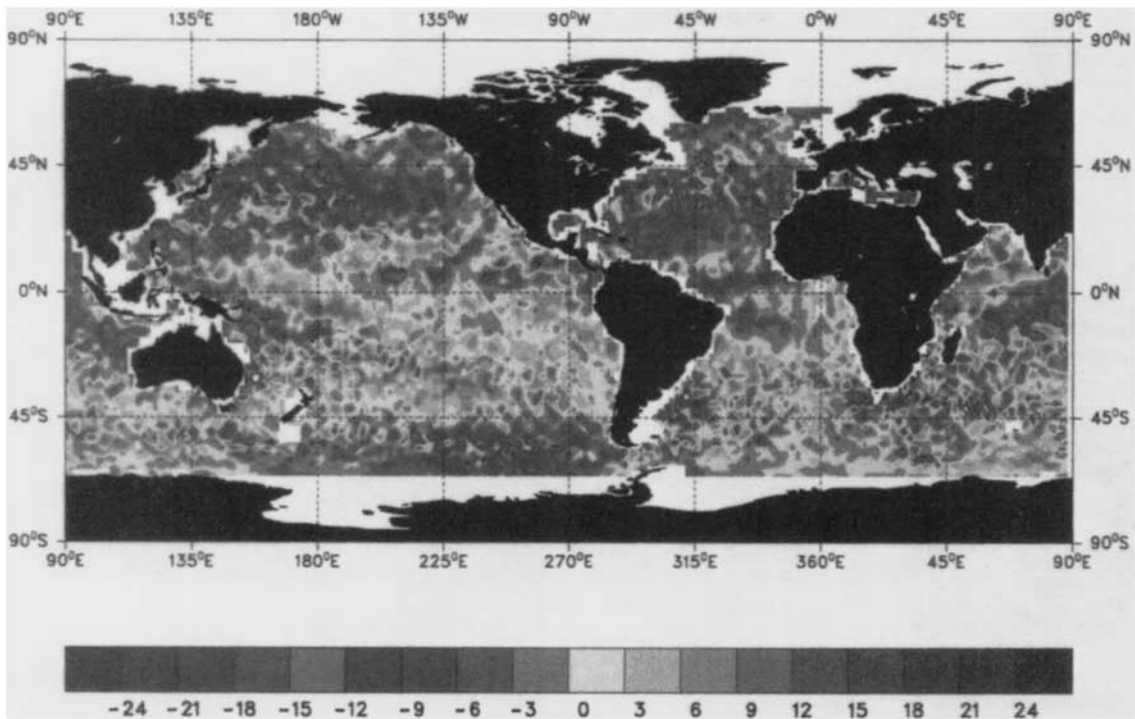


Figure 1-1b. Surface elevation of the ocean in cms during 10-20 March 1993 as seen by the TOPEX/POSEIDON altimeter satellite (after Stammer & Wunsch, 1994). The elevation is relative to a

2-year mean, has been averaged over 2° squares, and the 10-day “window” blurs the most rapidly changing features, therefore rendering the result somewhat simpler than a true instantaneous picture.

Nonetheless, both data and model confirm the essentially turbulent nature of the circulation. (A color reproduction of Figure 1-1b is shown in the color insert.)

same character in science, but unless recognized for what they are, they ultimately inhibit progress because it is not always clear to the student where the critical questions lie, what is really known, and what is merely assumed for convenience.

One of the great achievements of oceanography and fluid dynamics, in partnership with meteorology and other branches of geophysics, was the creation in the period following about 1950 of what is now known as *geophysical fluid dynamics*. This field created a dynamical and mathematical framework for discussion of the circulation, one that hardly existed prior to that time (see the remarks by Stommel, 1982). The elegance and rigor of this branch of fluid mechanics has permitted the growth of discussions of theories of the circulation which are interesting, useful, and even beautiful. But a side effect of this “applied mathematics of oceanography” has been to

create two, sometimes nearly independent and contradictory, views of the ocean: the ocean as observers understand it, and the ocean as the theoreticians describe it. Part of my motivation in writing this book has been to try to bring these two parts of the subject back into a closer relationship. I believe that we will progress most rapidly by continually calling attention to the real or apparent discrepancies between the theoretical picture and what we think observations imply.

Much of what we think we know about the ocean at any time is a direct consequence of the technology available at that point in history. One can usefully recall that in the Nansen bottle era, everyone knew that oceanic profiles of temperature and salinity were smooth functions of depth. The demonstration beginning in the middle 1960s with salinity/conductivity-temperature-depth profiling devices (STDs and CTDs) of what is now called *fine* and *microstructure* was initially greeted in many quarters with firm disbelief (the observations being attributed to faulty instruments). Until the advent of long time series of currents from drifting floats and moored current meters, everyone “knew” that nonsynoptic hydrography could be combined to produce a picture of the ocean circulation that was both steady and a true climatological average.

In general, the message is that as in all science, our understanding of the ocean comes through a distorting prism of our observational and theoretical technology; the student should maintain a very substantial degree of skepticism about almost anything said in textbooks about the ocean (beyond its being wet, salty, describable by the Navier-Stokes equations, and interesting), including this one.

Finally, oceanography is now struggling to become a true global science. The field is moving rapidly from what many now view with nostalgic regret as the romantic period of exploration—when “real” oceanographers were those who made all their own observations at sea and whose interpretations were limited mainly to drawing pictures. We are entering an era when the subject will necessarily become more like meteorology—with global data bases, obtained from observational networks run by large, impenetrable governmental agencies. The science will become less romantic, regrettably so (although like much of nostalgia, the romance of oceanography is often more apparent in the retelling than it was at the time). The excitement will still be there (it will be intellectual excitement), we will learn much more about what is going on, cherished myths will go on the trash heap, and new views of the ocean will surely emerge.

Some scientific fields, where there is little or no supporting theoretical framework, develop world views or paradigms through a process of scenario

development: Observations are made and a plausible story is told about how those data might be explained.¹ Subsequent observations may lead to elaborations of the original story, but a new or conflicting story is not told unless there is overwhelming evidence that the old one cannot be sustained. Fields that develop in this mode include geology and much of descriptive physical oceanography.

Until the development of geophysical fluid dynamics, physical oceanography did not have an adequate theoretical framework, and the scenario approach was the only one possible. Despite the existence of detailed understanding of fluid flows, much of physical oceanography has continued in the story-telling mode, as perusal of most any issue of a modern journal will show. The development of inverse and related methods—the focus of this book—is an attempt to shift the emphasis from the scenario-based methodology to one whose emphasis is more upon the quantitative (that is, numerical) testing of data against the known equations of fluid dynamics (and of chemistry and biology where available).² The mere use of numbers does not, however, embody the required shift: One often sees scenario-based discussions of the ocean circulation based on use of many millions of numbers acquired from computer models. It is the quantitative testing, and potential rejection, of models and scenarios using estimation methods that is required. It is only in this way that one can expect to see a convergence upon a consensus view of the ocean circulation, and an escape from the present circumstance in which dozens of apparently conflicting estimates of the circulation exist. As more data and understanding accumulate, uncertainty estimates should shrink, and prior estimates would be identified as either in contradiction with the new picture or consistent within the uncertainties.

Over the years, oceanography has benefited from being a junior partner of meteorology. Meteorology is much more mature, as a result of the technically easier problem of observing the atmosphere and the drive by governments to forecast the weather. There are many analogies between the two fluid systems (e.g., Charney & Flierl, 1981) and for which meteorological insight has proven most helpful to understanding the much more opaque,

¹ Some of these scenarios are reminiscent of Rudyard Kipling's *Just-So Stories*: "How the Leopard Got His Spots," "How the Camel Got His Hump," etc. We have "How the Ocean Circulates..."

² Recall the words of William Thomson, Lord Kelvin: "When you can measure what you are speaking about, and express it in numbers, you know something about it; but when you cannot measure it, when you cannot express it in numbers, your knowledge is of a meager and unsatisfactory kind: it may be the beginning of knowledge, but you have scarcely in your thoughts, advanced to the stage of *science*."

literally and figuratively, ocean. As one looks to the future, however, some divergences are apparent.

Much meteorological effort is directed at operational forecasting problems. It is not pejorative to characterize this work as engineering in part. If a forecasting technique works well, one should use it, whether or not it is fully understood. If the forecasts are going astray, the forecaster normally gets fast and vehement feedback from the public. Thus, what the meteorologist refers to as *assimilation*, the combination of model analysis with data for the purpose of making the best forecast, has tended to pay comparatively little attention to the uncertainty estimates of the process—if the forecasts work well, one knows it and that is enough. But the oceanographer has no such clientele (at the moment), the goals are directed more at understanding the system than forecasting it, and when I deal with assimilation methods, much more emphasis is placed upon sensitivity and uncertainty than a meteorologist might regard as necessary or even sensible. Here, the view is taken that because the error or uncertainty estimates describe what we don't know, they may be more important than the formal solution to a set of model equations. So some of this book addresses techniques that differ from normal meteorological practice.

The book is intended to be at a level accessible to first- or second-year graduate students with only the beginnings of a knowledge of ocean dynamics. I hope, too, that it will prove a useful guide to workers on the edge of physical oceanography, including biologists and chemists, as well as those entering from fields such as satellite altimetry and scatterometry, who seek some guidelines as to what is important and what is peripheral to this subject.

Beginning in the middle 1970s, it began to be realized that the classical *dynamic method* could be readily extended to produce absolute estimates of the ocean circulation. The thermal wind balance, coupled with such simple statements that mass and salt are conserved, largely sets the strength of the ocean circulation. The gross tilt of thermocline throughout the world ocean is one of the best known of all oceanic phenomena. It sets the clock that determines how fast the system is moving and what the ocean transports and exchanges with the atmosphere.

To some extent, this book is a tribute to the power of geostrophy—in the form of the thermal wind balance. The simple assertion that the vertical derivatives of the horizontal velocity are proportional to the horizontal derivatives of the density field is so familiar that the enormous quantitative power of the relationship tends to be submerged in discussions of measurements that are technically more interesting (e.g., of transient tracers).

But temperature and salinity are tracers, too; they are the easiest of all to measure; we have much better coverage of the world ocean for them than anything else; they are an intimate element of the global climate system; and their immediate relation to the density field means that they must be the central focus of any effort to understand the general circulation.

1.2 What Is an Inverse Problem?

What I mean by the title of this book, *The Ocean Circulation Inverse Problem*, is the problem of inferring the state of the ocean circulation, understanding it dynamically, and even perhaps forecasting it, through a quantitative combination of theory and observations. It may help the reader in what follows to understand why it is called an inverse problem and why the label is partially a misnomer, and to connect the ocean circulation problem to the many other such problems both in oceanography and in science as a whole. In particular, I wish to emphasize the difference between an inverse problem and what we are really discussing in this book—what are called inverse *methods*.

We digress slightly and explore some conventions of mathematics as applied to familiar differential systems, relying heavily on the reader’s experience. Consider a very familiar problem:

Solve

$$\nabla^2 \phi = \rho \tag{1.2.1}$$

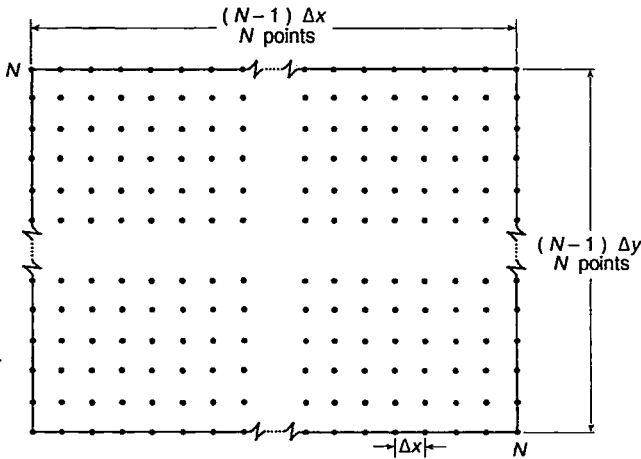
for ϕ , given ρ , in the domain $\mathbf{r} \in D$, subject to the boundary conditions $\phi = \phi_0$ on the boundary ∂D , where \mathbf{r} is a spatial coordinate.

This statement is the Dirichlet problem for the Laplace-Poisson equation, whose solution is well-behaved, unique, and stable to perturbations in the boundary data, ϕ_0 , and the source or forcing, ρ . Because it is a familiar boundary value problem, it is labeled a *forward* (or *direct*) problem.

Now consider a different version of the above: Solve (1.2.1) for ρ given ϕ in the domain D .

This latter problem is even easier to solve than the forward problem: merely differentiate ϕ twice to obtain the Laplacian, and ρ is obtained directly from (1.2.1). Because the problem as stated is inverse to the conventional forward one, it is labeled an *inverse problem*. It is inverse to a more familiar boundary value problem in the sense that the usual unknowns ϕ have been inverted or interchanged with (some of) the usual knowns ρ . Notice

Figure 1–2. Simple square, homogeneous grid used for discretizing the Laplacian, reducing the differential equation to a set of simultaneous equations.



that both problems, as posed, are well-behaved and produce uniquely determined answers (ruling out mathematical pathologies in any of ρ , ϕ_0 , ∂D , or ϕ). This is not the only inverse problem that could be set; we could, for example, demand computation of the boundary conditions, ϕ_0 , from given information about some or all of ϕ , ρ .

Because this book is based upon discrete methods, write the Laplace-Poisson equation in finite difference form for two Cartesian dimensions:

$$\phi_{i+1,j} - 2\phi_{i,j} + \phi_{i-1,j} + \phi_{i,j+1} - 2\phi_{i,j} + \phi_{i,j-1} = \rho_{ij}, \quad i, j \in D. \tag{1.2.2}$$

To make the bookkeeping as simple as possible, suppose the domain D is the square $N \times N$ grid displayed in Figure 1–2, so that ∂D is the four line segments shown. There are $(N - 2) \times (N - 2)$ interior grid points, and Equations (1.2.2) are then $(N - 2) \times (N - 2)$ equations in N^2 of the ϕ_{ij} . If this is the forward problem with ρ_{ij} specified, there are fewer equations than unknowns. But if we append to (1.2.2) the set of boundary conditions:

$$\phi_{ij} = \phi_{ij}^0, \quad i, j \in \partial D, \tag{1.2.3}$$

there are precisely $4N - 4$ of these conditions, and thus the combined set (1.2.2) plus (1.2.3), which we write as

$$\mathbf{A}_1 \phi = \mathbf{d}_1, \tag{1.2.4}$$

is a set of $M = N^2$ equations in $M = N^2$ unknowns, with

$$\phi = \begin{bmatrix} \phi_{11} \\ \phi_{12} \\ \vdots \\ \vdots \\ \phi_{NN} \end{bmatrix}, \quad \mathbf{d}_1 = \begin{bmatrix} \rho_{11} \\ \rho_{12} \\ \vdots \\ \vdots \\ \rho_{N-2,N-2} \\ \phi_{11}^0 \\ \vdots \\ \phi_{ij}^0 \end{bmatrix}.$$

The nice properties of the Dirichlet problem can be deduced from the well-behaved character of the matrix \mathbf{A}_1 . Thus the forward problem corresponds directly with the solution of an ordinary set of simultaneous algebraic equations (Lanczos, 1961, has a much fuller discussion of this correspondence).

A complementary inverse problem says: “Using (1.2.4), compute ρ_{ij} and the boundary conditions, given ϕ_{ij} ,” an even simpler computation—it involves just multiplying the known ϕ by the known matrix \mathbf{A}_1 . The problem can be written formally as

$$\mathbf{A}_2 \mathbf{d}_1 = \mathbf{d}_2,$$

where \mathbf{A}_2 is the identity, and $\mathbf{d}_2 = \mathbf{A}_1 \phi$. Presumably all this is obvious. But now let us make one small change in the forward problem, changing it to the Neumann problem:

Solve

$$\nabla^2 \phi = \rho \tag{1.2.5}$$

for ϕ , given ρ , in the domain $\mathbf{r} \in D$ subject to the boundary conditions $\partial \phi / \partial \mathbf{m} = \phi'_0$ on the boundary ∂D , where \mathbf{r} is a spatial coordinate and \mathbf{m} is the normal to the boundary.

This new problem is another classical, much analyzed forward problem. It is, however, well-known that the solution to (1.2.5) with these new boundary conditions is indeterminate up to an additive constant. This indeterminacy is clear in the discrete form: Equations (1.2.3) are now replaced by

$$\phi_{i+1,j} - \phi_{i,j} = \phi'_{ij}, \quad i, j \in \partial D' \tag{1.2.6}$$

etc., where $\partial D'$ represents the set of boundary indices necessary to compute the local normal derivative. There is a new combined set:

$$\mathbf{A}_3 \phi = \mathbf{d}_3. \tag{1.2.7}$$

Because only *differences* of the ϕ_{ij} are specified, there is no information

concerning the mean value of ϕ . When we obtain some proper machinery in Chapter 3, we will be able to demonstrate that even though (1.2.7) appears to be M equations in M unknowns, in fact only $M - 1$ of the equations are independent, and thus the Neumann problem is an underdetermined one. This property of the Neumann problem is well-known, and there are many ways of handling it, either in the continuous or discrete forms. In the discrete form, a simple way is to add one equation setting the value at any point to zero (or anything else).

Notice however, that the inverse problem remains unchanged, well-posed, and unique. Although somewhat trivial in form, we can write it as a set of simultaneous linear equations, by mere rearrangement of (1.2.7),

$$\mathbf{A}_4 \mathbf{d}_3 = \mathbf{d}_4 \tag{1.2.8}$$

where \mathbf{A}_4 is the identity matrix and $\mathbf{d}_4 = \mathbf{A}_3 \phi$.

These are examples in which a forward problem is badly posed in the sense of missing a piece of information necessary to determine the solution uniquely, while the inverse problem is fully posed. This point is labored a bit, because later we will encounter some inverse problems that are also missing some of the relevant information. Among the techniques for solving problems which are underdetermined are a class sometimes known as *inverse methods*. But one must carefully distinguish between inverse *problems* and the *methods* available for solving undetermined systems of equations, whatever the origin of the problem—whether forward or inverse.

Many examples of inverse problems are mathematically and computationally well behaved. One elegant example is Abel’s problem (Aki & Richards, 1980): An observer stands on the floor of a symmetric valley and wishes to determine the shape of the valley. For reasons we need not inquire into, his method is to set a ball in motion, with initial speed s_0 ; he then measures the time it takes for the ball to return to him, yielding a set of travel times $t(s_0)$. Abel (1826) showed that the shape of the valley, $f(x)$, could be obtained as the solution to the integral equation (nondimensionalized):

$$t(s_0) = \int_0^{s_0} \frac{f(\xi)}{\sqrt{s_0 - \xi}} d\xi \tag{1.2.9}$$

(an Abel integral equation), with solution

$$f(\xi) = -\frac{1}{\pi} \frac{d}{d\xi} \int_\xi^a \frac{t(x')}{\sqrt{x' - \xi}} dx', \tag{1.2.10}$$

a stable, well-behaved solution. So the Abel problem is an example of an in-