



Marine Environmental Data Analysis

Hailong Liu

April 14 2022

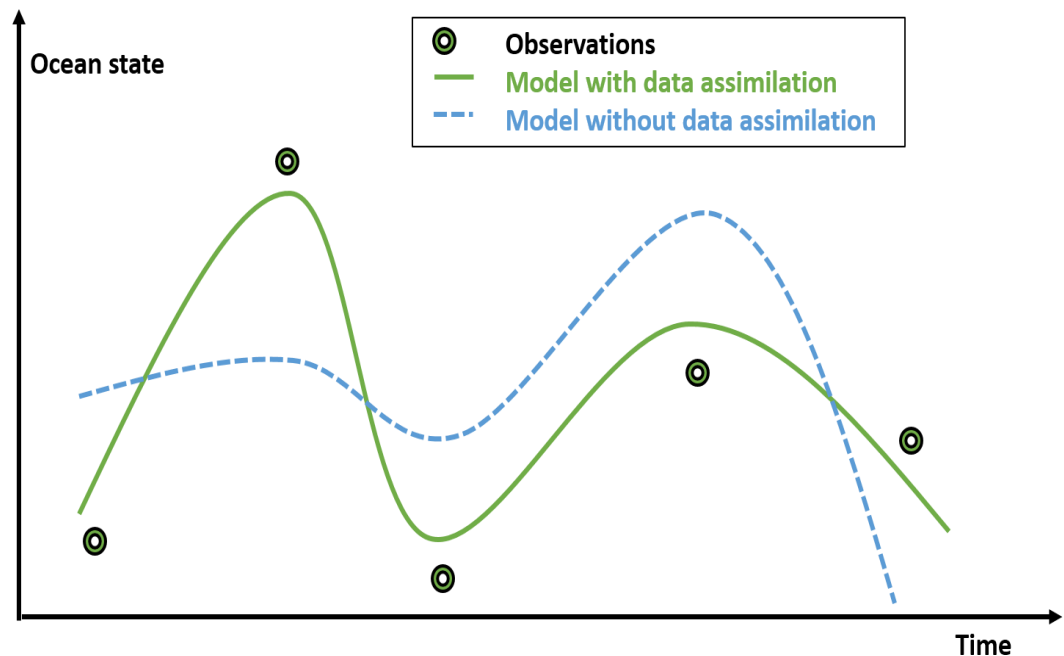


Data Sources

- Observations
- Numerical Model Simulations
- Reanalysis?

Reanalysis

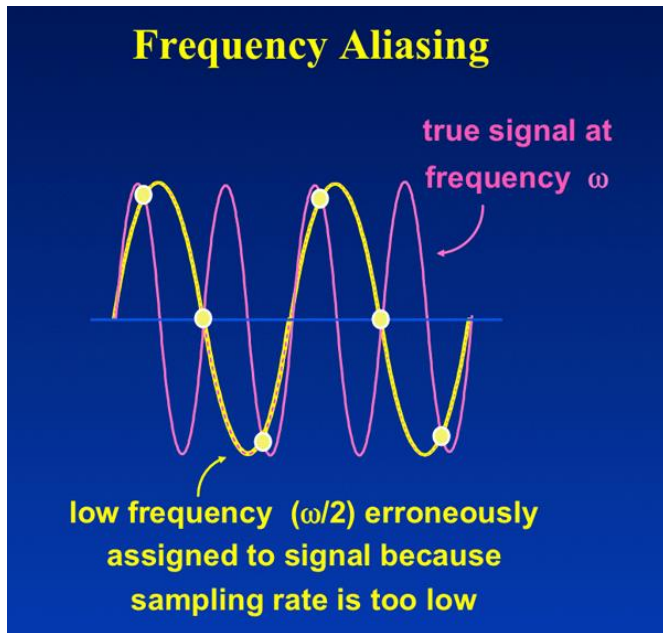
- Model + Observation



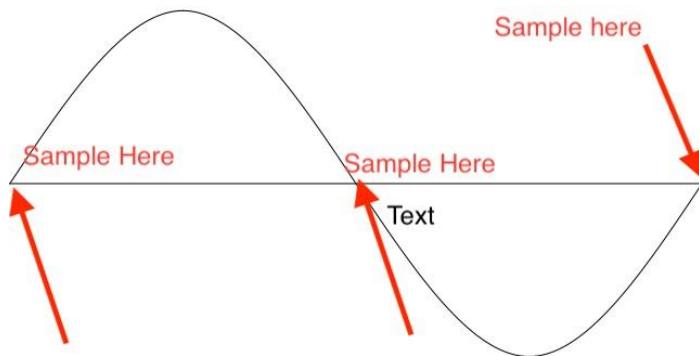
Atmosphere: NCEP/NCAR Reanalysis ...

Ocean: Simple Ocean Data Assimilation Ocean/Sea Ice Reanalysis (SODA) ...

Chapter 1 Data Acquisition and Recording



Sine with period of T
Sampling at T/2



The important aspect to keep in mind is that, for a given sampling interval Δt , the highest frequency we can hope to resolve is the *Nyquist* (or *folding*) frequency, f_N , defined as

$$f_N = 1/2\Delta t \quad (1.1)$$

total record length $T = N\Delta t$ obtained for N data samples that: (1) determines the lowest frequency (*the fundamental frequency*)

$$f_0 = 1/(N\Delta t) = 1/T \quad (1.2)$$

($\Delta f = 0.00305$ cph). Note that since f_N is the highest frequency we can measure and f_0 is the limit of our frequency resolution, then

$$f_N/f_0 = (1/2\Delta t)/(1/N\Delta t) = N/2 \quad (1.3)$$

is the maximum number of Fourier components we can hope to estimate in any analysis.



Chapter 1 Data Acquisition and Recording

Temperature (in-situ and potential)

in the ocean. To a high degree of approximation, the potential temperature defined as $\theta(p)$, or T_θ , is given in terms of the measured in situ temperature $T(p)$ as $\theta(p) = T(p) - F(R)$, where $F(R) \approx 0.1^\circ\text{C}/\text{km}$ is a function of the adiabatic temperature gradient R . The results can have important

Unit:

$$^\circ\text{F (Fahrenheit)} = 9/5^\circ \times \text{Celsius} + 32$$

$$^\circ\text{C (Celsius)} = (^\circ\text{F} - 32) \times 5/9$$

$$\text{K (kelvin)} = ^\circ\text{C} + 273.15 \text{ (0 K = absolute zero)}$$



Chapter 1 Data Acquisition and Recording

Salinity

samples. In its report (Forch et al., 1902), the commission recommended that salinity be defined as follows: “The total amount of solid material in grams contained in 1 kg of seawater when all the carbonate has been converted to oxide, all the bromine and iodine replaced by chlorine, and all the organic material oxidized.”

Unit

bles referred to as Knudsen’s tables. The symbol ‰ indicates “parts per thousand” (ppt) in analogy to percent (%) which is parts per hundred. In the more modern Practical Salinity Scale, salinity is a unitless quantity but, because many oceanographers are uncomfortable with unitless values, salinity is sometimes written as “psu” for *practical salinity units*.



Chapter 1 Data Acquisition and Recording

Absolute Salinity

...In 2010 a new standard for the properties of seawater called the *thermodynamic equation of seawater 2010* (TEOS-10) was introduced, advocating absolute salinity as a replacement for practical salinity, and [conservative temperature](#) as a replacement for [potential temperature](#).

Absolute salinities on this scale are expressed as a mass fraction, in grams per kilogram of solution. Salinities on this scale are determined by combining electrical conductivity measurements with other information that can account for regional changes in the composition of seawater...



Density and Potential Density

The use of potential temperature in the calculation of density leads to the definition of potential density, $\rho_\theta = \rho(S, \theta, 0)$ in kg/m^3 , as the value



Raw data

- Calibration
- Errors Check
 1. large accidental errors/spikes that result from equipment failure > subjective evaluation
 2. small random errors or noise > statistical methods
- Interpolation (most analysis methods require data values that are regularly spaced in time or space)
- Data presentation (visual display)

Vertical Profiles

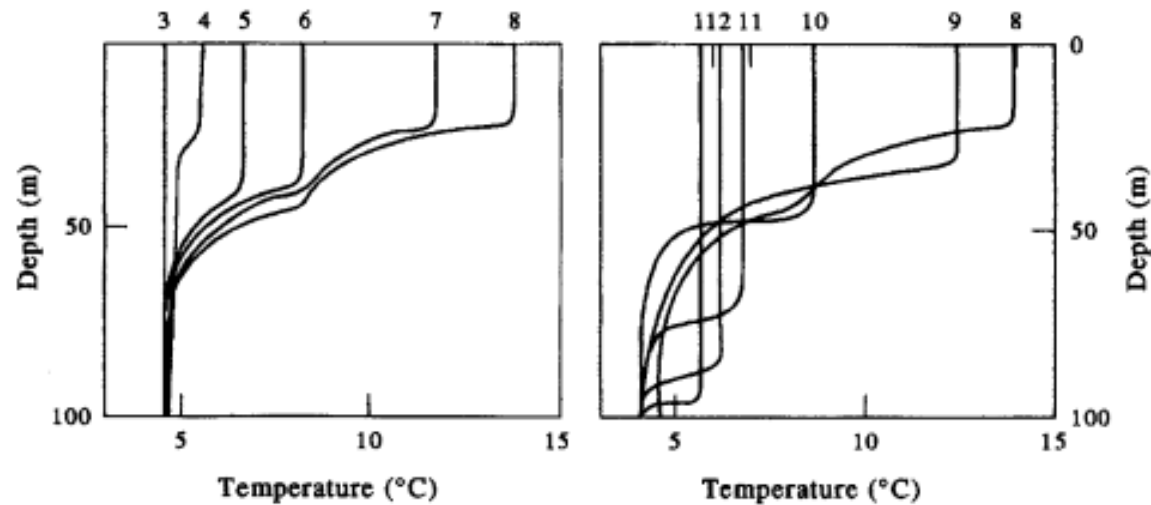
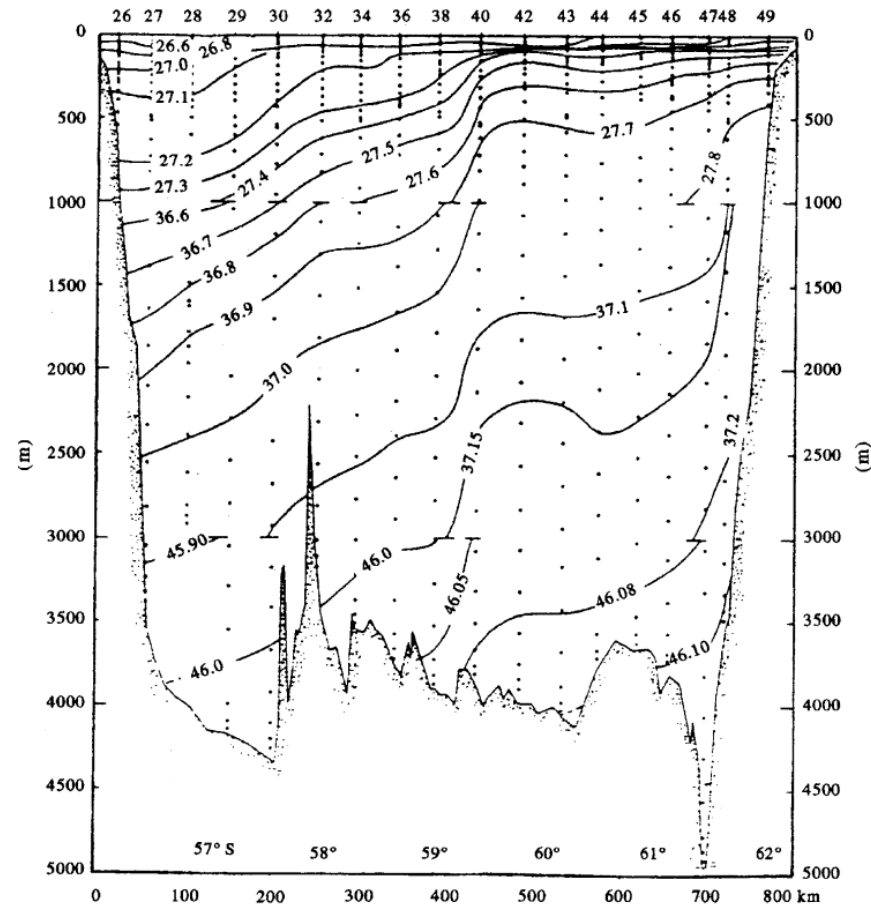


FIGURE 2.6 Time series of monthly mean profiles of upper ocean temperature at Ocean Weather Station "P", northeast Pacific (50° N, 145° W). Numbers denote the months of the year. (From Pickard and Emery (1992).)

Vertical Sections



JRE 2.7 Cross-section of density (σ) (kg/m^3) across Drake Passage in 1976. (From Nowlin et al. (1

Time Series

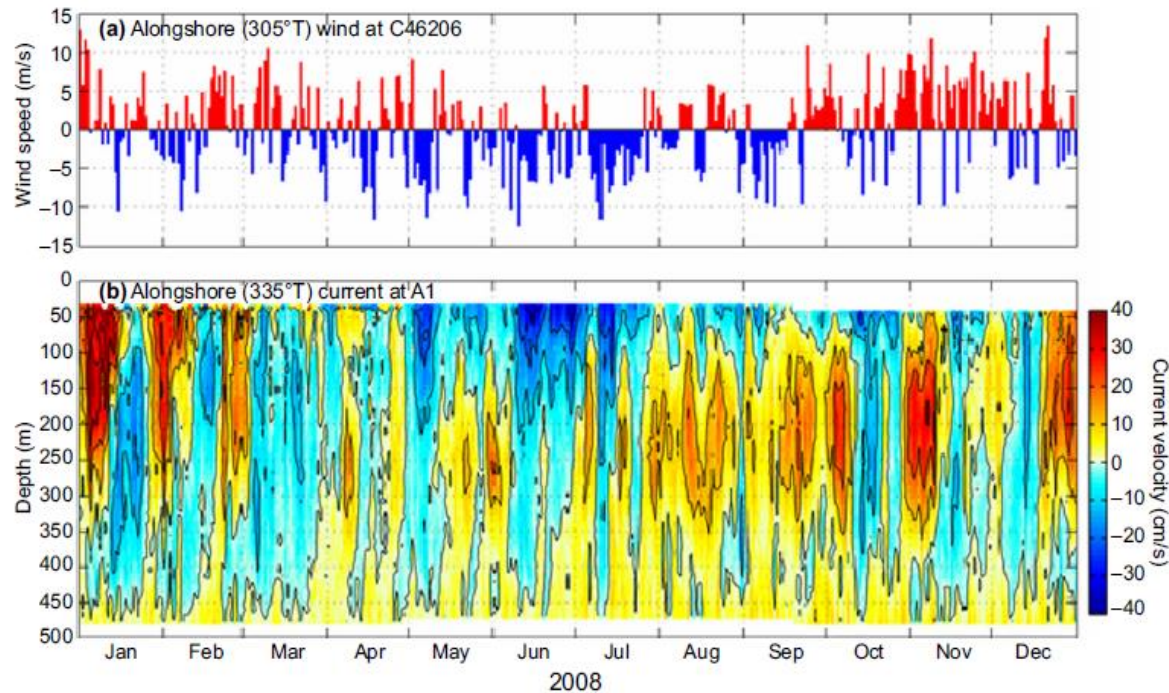
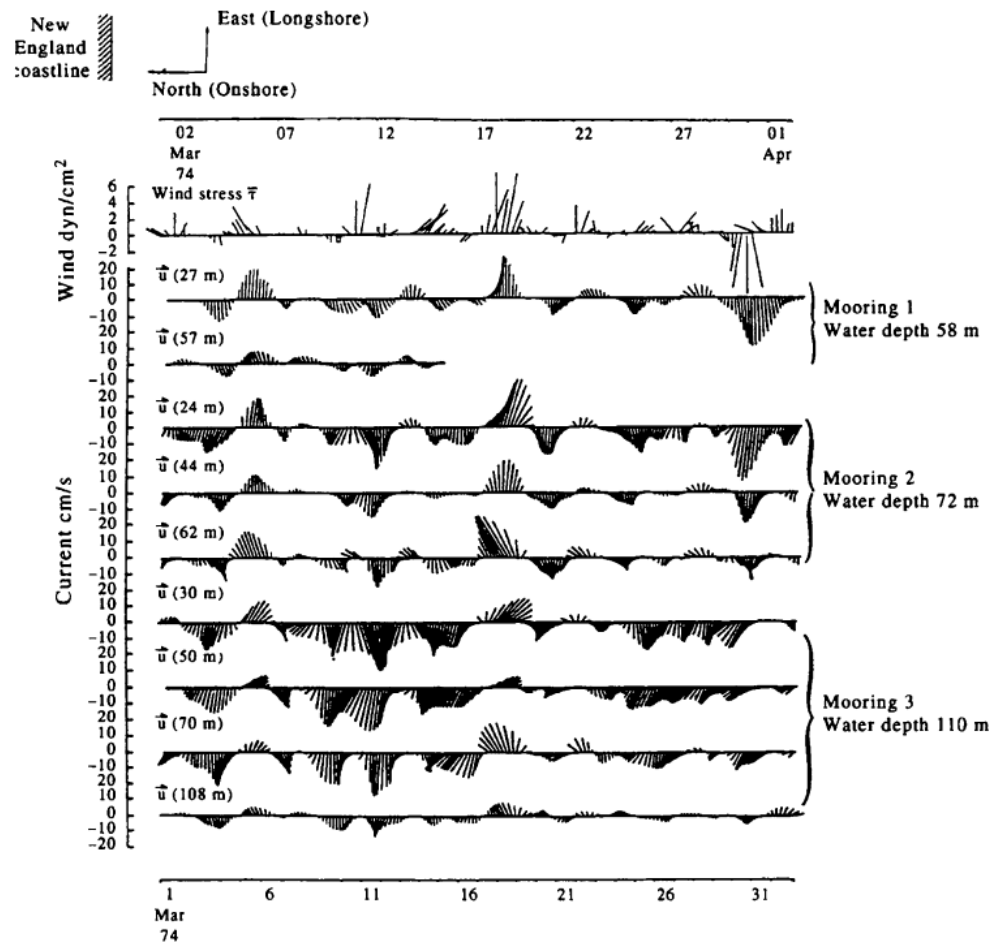


FIGURE 2.18 Time series (January 1–December 31, 2008) of the daily mean alongshore component of (a) wind stress at meteorological buoy C46206 and (b) current velocity at mooring site A1. The locations are within 20 km of one another off southwest Vancouver Island, British Columbia. The positive alongshore directions (305 and 335 °T, respectively) are given in “degrees True” compass bearing. Currents are measured every 4 m through the water column from an upward-looking 75-kHz ADCP moored in 500 m of water on the continental slope. Winds are measured about 20 km away on the continental slope. (*Data and analysis courtesy of Maxim Krassovski, Steve Mihály, Tamás Juhász and David Spear.*)

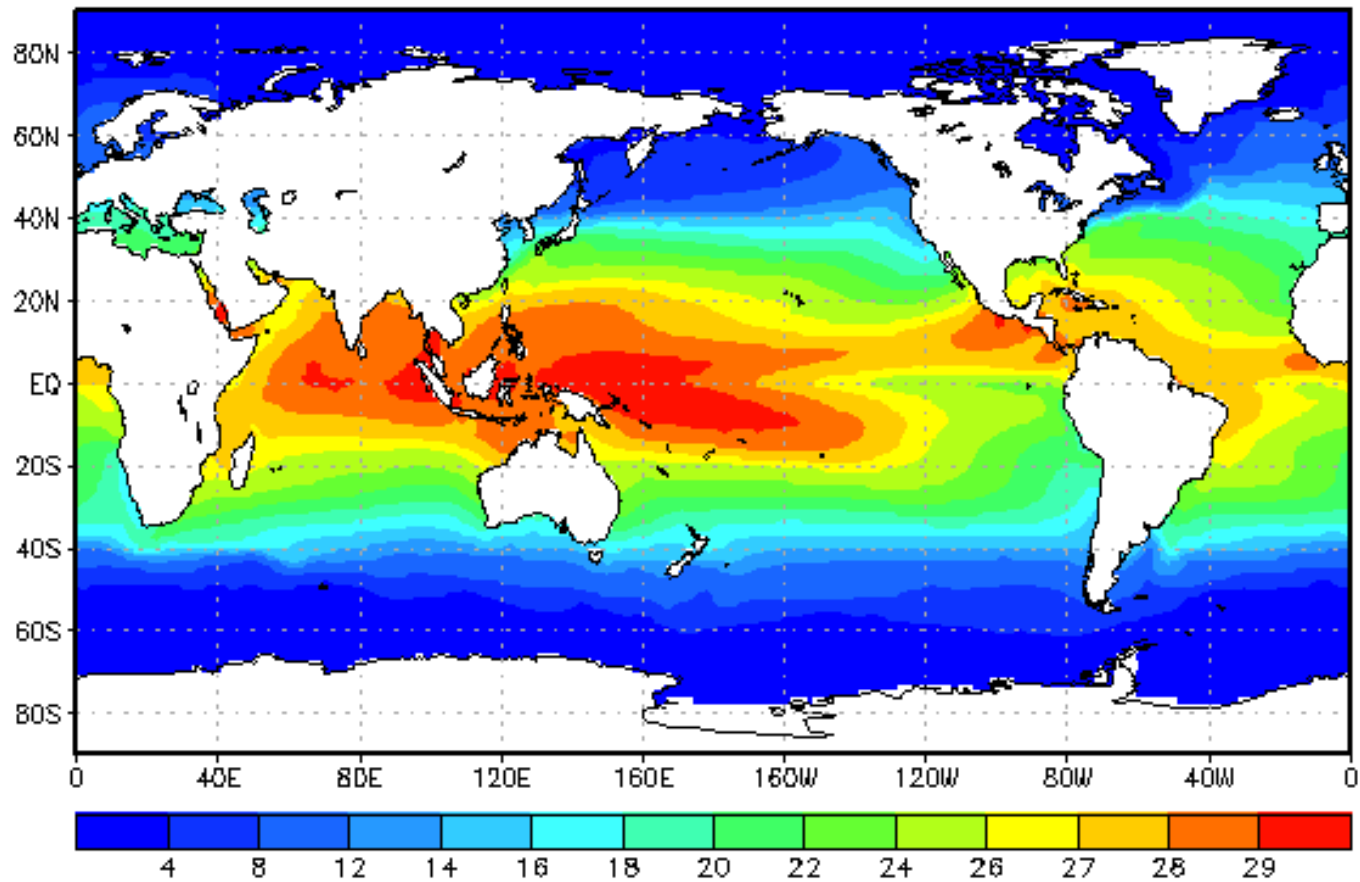
Chapter 2 Data Processing and Presentation

Vector Plots



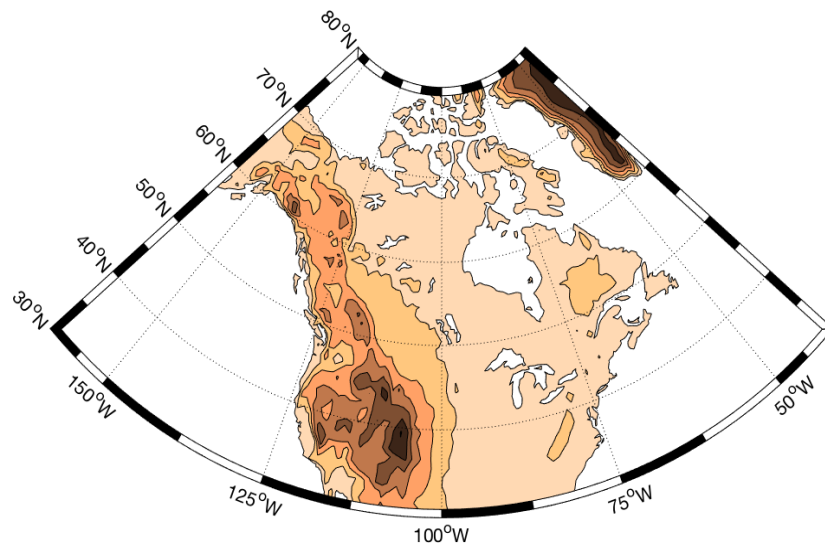
› Vector (stick) plots of low-pass filtered wind stress and subtidal currents at different depths moored off the United States about 100 km west of Nantucket Shoals. East (up) is alongshore and north (right) is onshore. The current meter depth (m). (Figure 7.11 from Beardsley and Boicourt (1981).)

Horizontal Maps

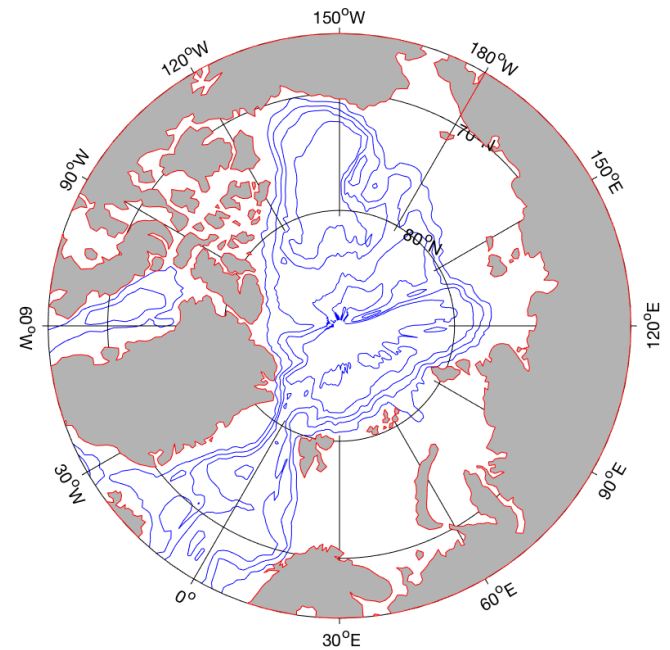


Map Projections

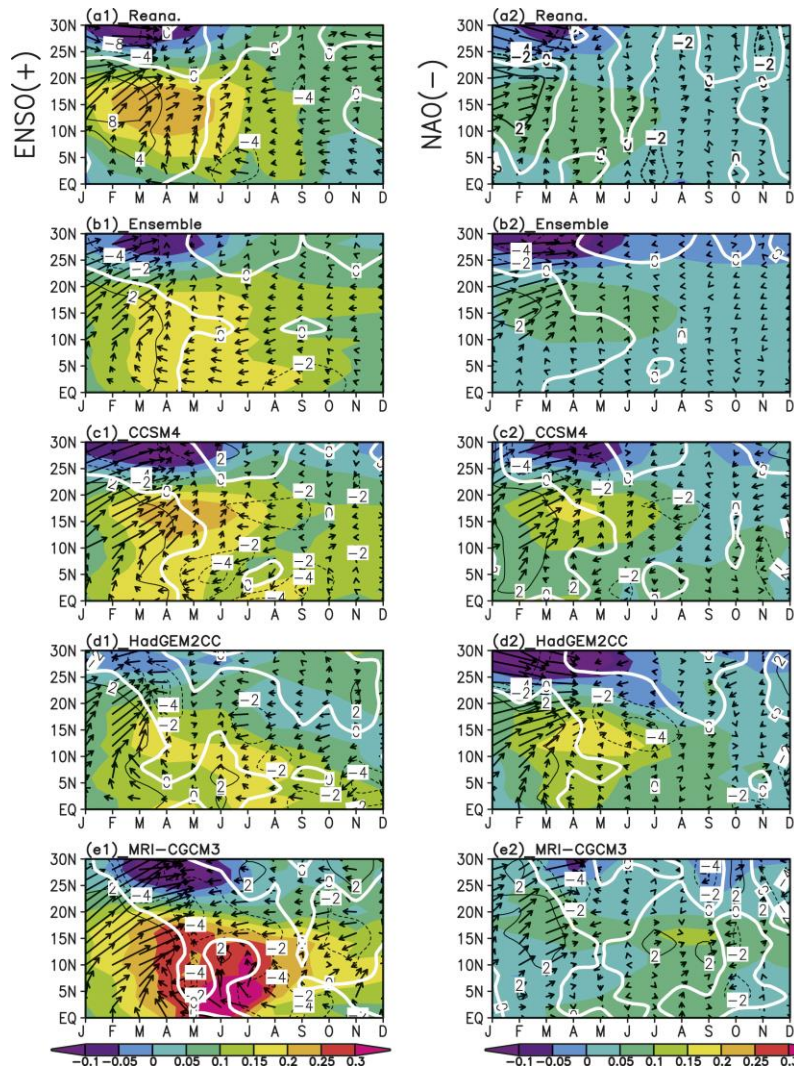
Lambert Conformal Conic Projection



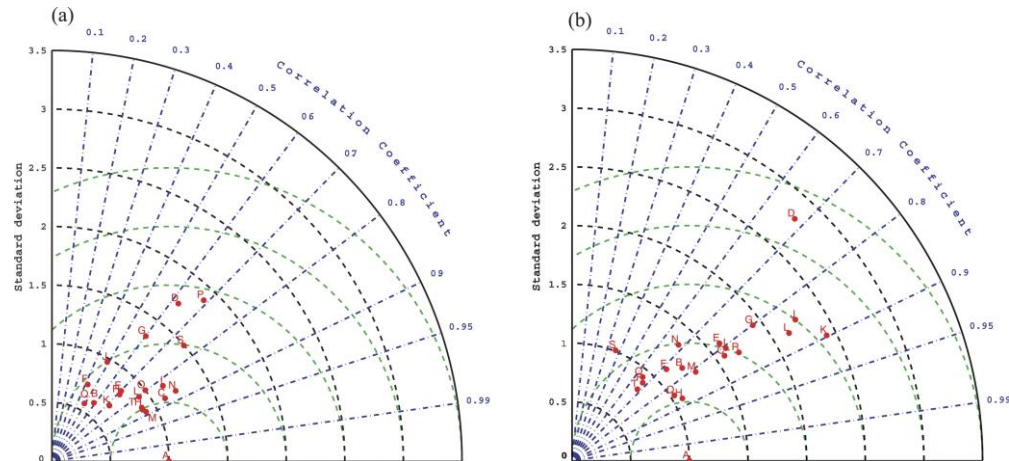
Stereographic Projection



Chapter 2 Data Processing and Presentation



Taylor Diagrams





Three reasons why a statistical approach is practical

1. The ocean is complex, requiring the specification of many variables, in fact, many more than can be observed or specified.
2. The ocean is nonlinear, so that groups of variables cannot be studied in isolation.
3. Ocean observations are not controlled; all variables change as the system evolves.



Chapter 3 Statistical Methods and Error Handling

Basic Notion

The goal of data analysis: discover relations and features in a dataset.

1.1 Expectation and mean

1.1.1 Continuous variables

For a random variables x which takes on continuous values over a domain Ω , the expectation is given by an integral,

$$E(x) = \int_{\Omega} xp(x)dx \quad (1.1)$$

where $p(x)$ is the probability density function. For any function $f(x)$, the expectation is

$$E(x) = \int_{\Omega} f(x)p(x)dx \quad (1.2)$$

Example: For a normal distribution: $x \sim N(0,1)$, the expectation of x

$$E(x) = \int_{-\infty}^{\infty} xp(x)dx = \int_{-\infty}^{\infty} x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = \int_0^{\infty} x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx + \int_{-\infty}^0 x \frac{1}{\sqrt{2\pi}} e^{-\frac{x^2}{2}} dx = 0$$

1.1.2 Discrete variables

Let x be random variable which takes on discrete. For example, x can be the outcome of a die cast, where the possible values are $x_i = i$, with $i=1,2,\dots,6$.

The expectation or expected value of x_i from a population is given by

$$E[x] = \sum_i x_i p_i \quad (1.3)$$

Where p_i is the probability of x_i occurring. If the die is fair, p_i is $1/6$ for all i ,

So $E[x] = \sum_i x_i p_i = (1+2+3+4+5+6)/6=3.5$. We also write

$$E[x] = \mu_x$$



Chapter 3 Statistical Meth

with μ_x denoting the mean of x for the population.

Similarly, for any discrete function $f(x)$, its expectation is

$$E[f(x)] = \sum_i f(x_i) p_i \quad (1.4)$$

The expectation of a sum of random variables satisfies

$$E[ax + by + c] = aE(x) + bE(y) + c, \quad (1.5)$$

where x and y are random variables, and a, b, c , are constants.

In practice, one can only sample N measurements of $x(x_1, \dots, x_N)$ from the population. The sample mean \bar{x} or $\langle x \rangle$ is calculated as

$$\bar{x} \equiv \langle x \rangle = \frac{1}{N} \sum_{i=1}^N x_i \quad (1.6)$$

which is in general different from the population mean μ_x . As the same size increases, the same mean approaches the population mean.

The function of “mean” in MATLAB is `mean(x)`

1.2 Variance and covariance

Fluctuations about the mean value is commonly characterized by the variance of the population,

$$\text{Var}(x) \equiv E[(x - \mu_x)^2] = E(x^2 - 2x\mu_x + \mu_x^2) = E[x^2] - \mu_x^2 \quad (1.7)$$

where (1.5) has been invoked. The standard deviation s is the positive square root of the population variance, i.e.,

$$s^2 = \text{Var}(x) \quad (1.8)$$

The sample standard deviation $\sigma^2 = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})^2$ (1.9)



Chapter 3 Statistical Methods

As the sample size increases, the sample variance approaches the population variance. For large N , distinction is often not made between having $N-1$ or N in the denominator of (1.9).

Normalization: The goal is to make different variables to be measured and compared in the same scale, i.e., an effective way to remove the influence of units.

For example, we would like to compare two very different variables, e.g., sea surface temperature and fish population. Simply, one can't even draw their variations in a plot due to different units. So, one usually standardizes the variables before making the comparison. The standardized variable

$$\mathbf{x}_s = (\mathbf{x} - \bar{\mathbf{x}}) / \sigma \quad (1.10)$$

is obtained from the original variables by subtracting the sample mean and dividing by the sample standard deviation. The standardized variable is also called the normalized variable or the standardized anomaly (where anomaly means the deviation from the mean value).

For two random variables x and y , with mean μ_x and μ_y respectively, their covariance is given by

$$\text{Cov}(x, y) = E[(x - \mu_x)(y - \mu_y)] \quad (1.11)$$

The variance is simply a special case of the covariance, with

$$\text{Var}(x) = \text{Cov}(x, x). \quad (1.12)$$

The sample covariance is computed as

$$\text{Cov}(x, y) = \frac{1}{N-1} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y}). \quad (1.13)$$

The function of variance, standard deviation and covariance in MATLAB is `var(x)`, `std(x)`, `cov(x,y)`.

1.3 Skewness:

Skewness is a measure of symmetry or more precisely, the lack of symmetry.

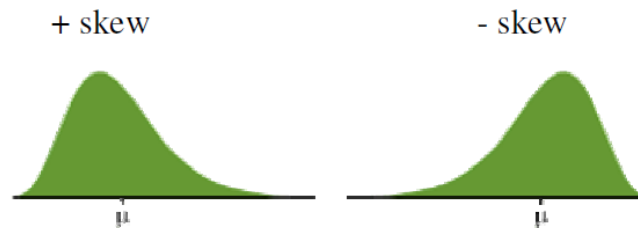


$$skew(x) = \frac{E[(x - \mu_x)^3]}{\sigma_x^3} \quad (1.14)$$

where μ_x and σ_x are the mean and standard deviation of x . As one might expect, the formula takes on a positive value if x is positively skewed and a negative value if x is negatively skewed

A distribution, or data set, is symmetric if it looks the same to the left and right of the center point, i.e., skew = 0. If the left **tail** (tail at small end of the the distribution) is more pronounced than the right tail (tail at the large end of the distribution), the skew is **negative**. If the reverse is true, the skew is **positive**. Distributions with positive skew are sometimes called "skewed to the right" whereas distributions with negative skew are called "skewed to the left."

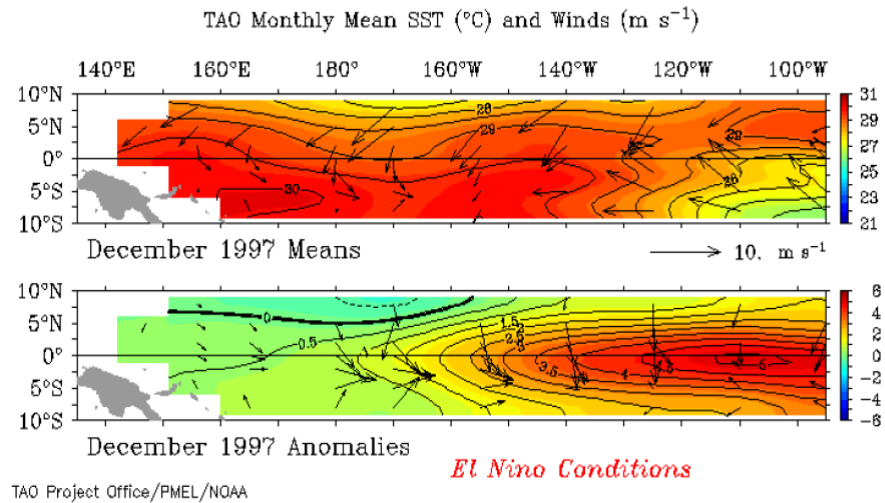
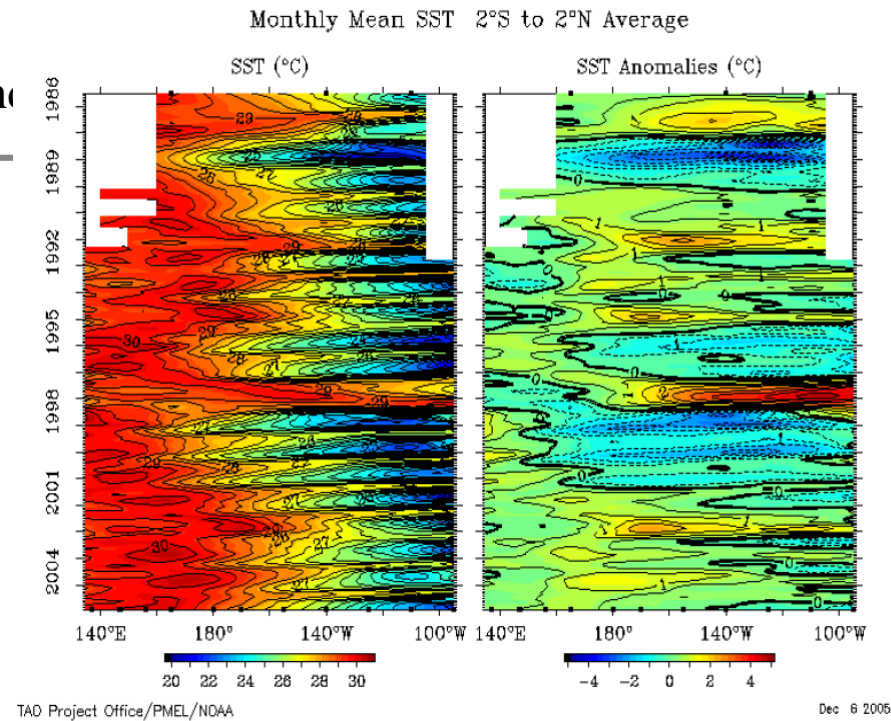
Positive vs. Negative Skewness



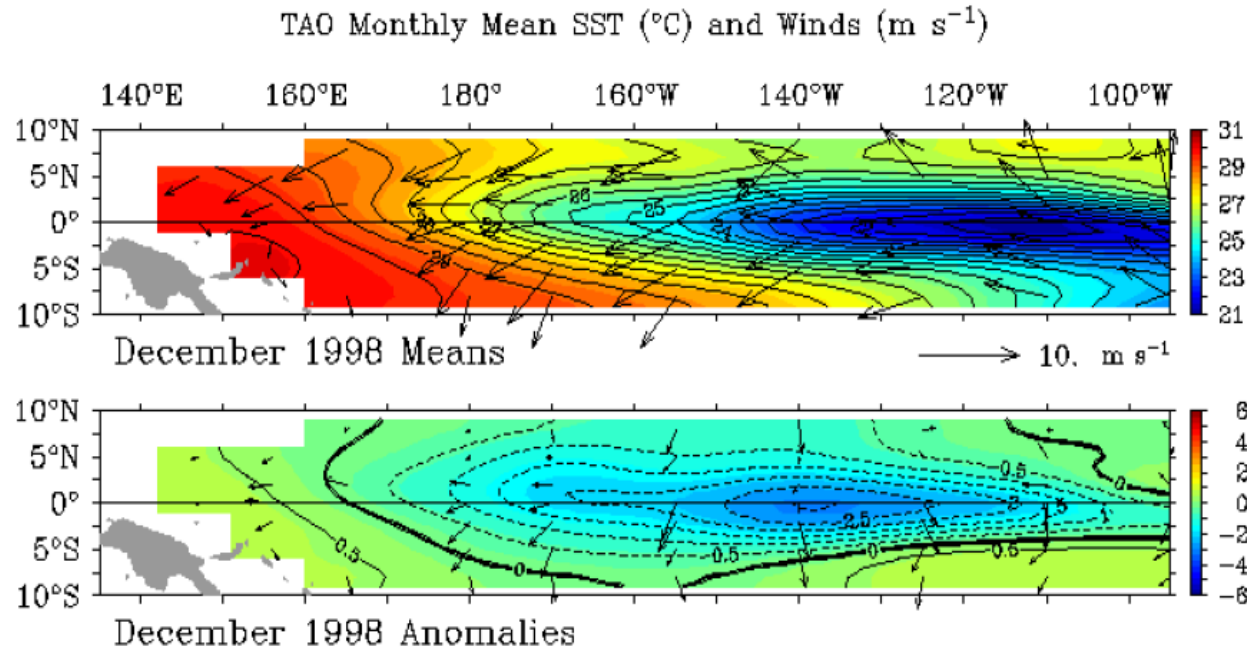
These graphs illustrate the notion of skewness. Both PDFs have the same expectation and variance. The one on the left is positively skewed. The one on the right is negatively skewed.

1.4 Examples of using these concepts to analysis practical problems.

Examples



Chapter 3



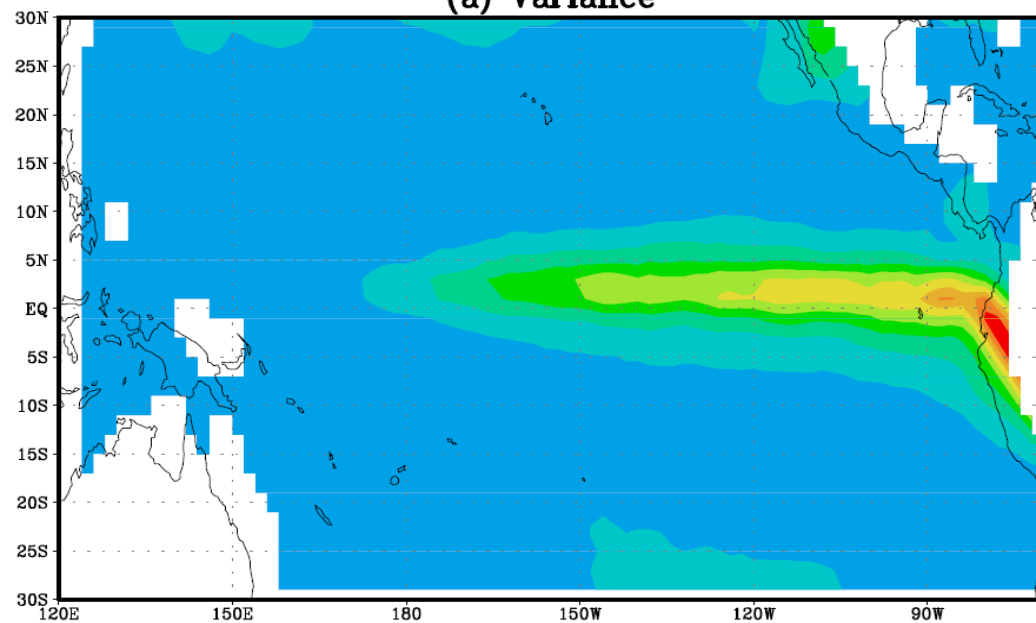
La Nina Conditions

TAO Project Office/PMEL/NOAA

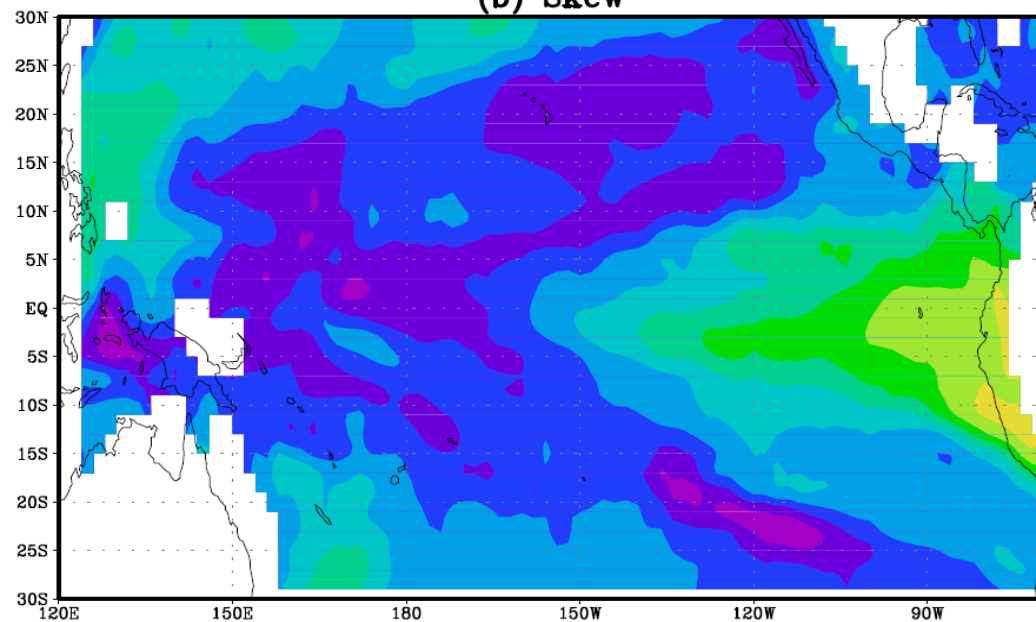
Shown above are two cases: El Nino and La Nina. Comparing the two cases reveals: (1) the largest variations occur at the eastern Pacific ocean; (2) El Nino and La Nina are asymmetric.

These two features could be explained by the below figure. As can be seen, the largest magnitudes of variances appear over the eastern Pacific, coinciding with the above figures. Similarly, the large positive skewness occupies the equatorial eastern boundary, and small negative skewness appears in the west, indicating the asymmetry shown between El Nino and La Nina.

(a) Variance



(b) Skew



Today class is over

