

# Multivariate Analysis of the North Atlantic Ocean - the problem of missing values

Peter Challenor  
Richenda Houseago-Stokes

## Sea Surface Temperature (SST)

---

- Meteorologists and oceanographers are interested in the interaction between the atmosphere and ocean
- The ocean influences the atmosphere through the sea surface temperature
- (the atmosphere's influence on the ocean is more complex)
- SST is measured from ships, buoys and **satellites**

## Sea surface temperature

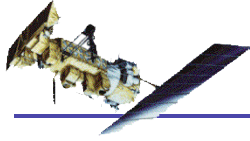
---

- From satellite SST we can identify and monitor surface disturbances that cross entire ocean basins, track ocean eddies and map ocean fronts
- It can also reveal striking features such as 'storms' in the upper ocean, known as eddies.
  - These are typically ~100 km wide and carry large amounts of energy around the globe. They play an important role in ocean circulation and climate.
- Space-borne infra-red sensors estimate SST by measuring heat radiation from the ocean surface.

## SST

---

- SST is measured using a infrared radiometer
- Spectral bands used are near the peak of surface emission - the peak ones aren't used due to atmospheric effects
- It is measured by:
  - taking the intensity of radiation at top of atmosphere
  - removing the atmospheric contribution
  - Resulting in the brightness temperature at the surface
- Brightness temp is approximately equal to the SST
- Typical SST products:
  - NOAA MCSST global
  - ATSR ASST



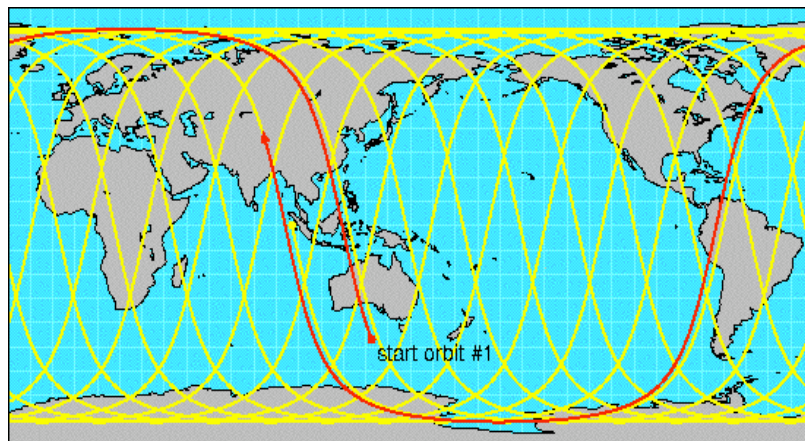
## AVHRR

- The AVHRR instruments are onboard the NOAA-7, -9, -11 and -14 satellites
- The radiometer uses 5 detectors that collect different bands of radiation wavelengths

0.63 $\mu\text{m}$	Daytime cloud and surface mapping
0.86 $\mu\text{m}$	Land-water boundaries
3.74 $\mu\text{m}$	Night cloud mapping, sea surface
10.80 $\mu\text{m}$	Night cloud mapping, sea surface
11.50 $\mu\text{m}$	Sea surface temperature

- Data from both the ascending pass (daytime) and descending pass (nighttime) are available with a spatial resolution of about 5.7 pixels per degree of longitude and latitude.

## Polar orbit track



Source: NOAA

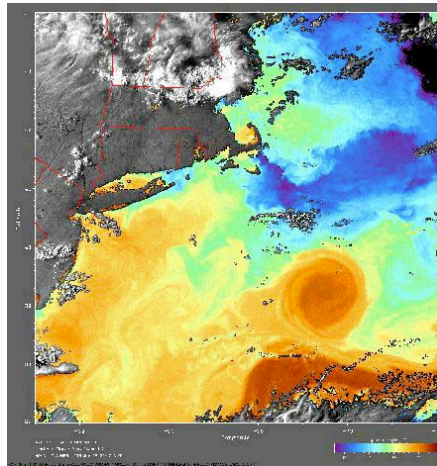
## Cloud

---

- Infrared radiation does not penetrate cloud
- Therefore where it is cloudy we get no data
- This means we have intermittent data

## SST image

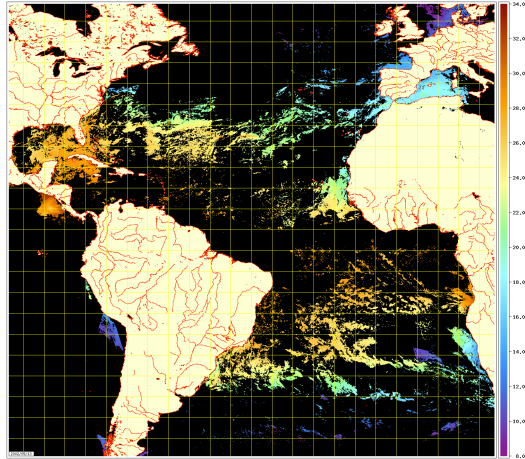
---



West Atlantic 15th June 1996 (NOAA)

## Cloud contamination

---



AVHRR image 13th May 2002 (NOAA)

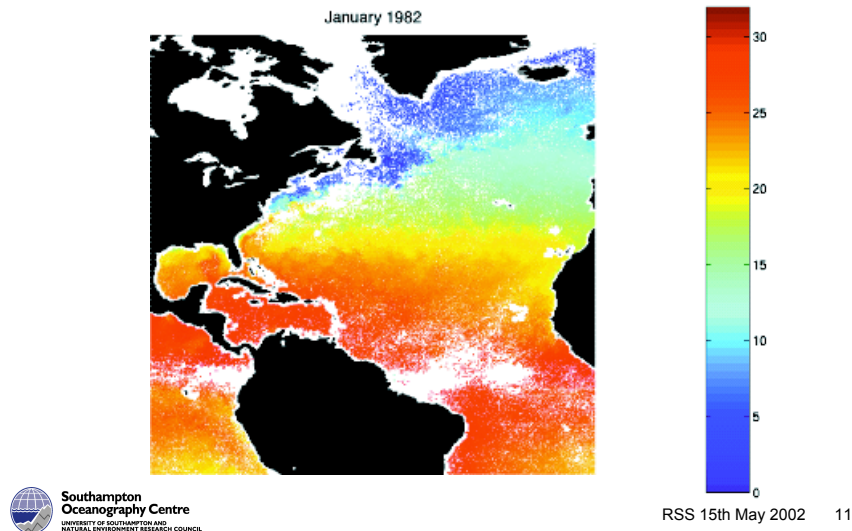
## MCSST Data

---

- Weekly averaged MCSST data are available from 11 November 1981 to 7 February 2001
- Data from January 1984 to December 2000 are being used - problem due to aerosols from the eruption of El Chichon
- Daytime and nighttime data is available.
- Areas, particularly north of 60°N, have few values due to extensive cloud cover and sea ice.
- MCSST data are produced at 18km resolution we average to 1°

## North Atlantic SST time series

---



## PCA with missing values

---

- Conventional principal component analysis cannot be done if there are missing values
- Two possible approaches:
  - Interpolate the data before doing the analysis
  - Use a modified version of PCA that can cope with missing values

## Interpolating the data

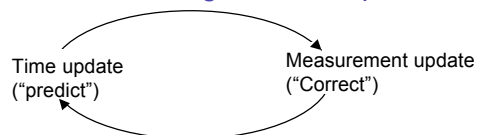
---

- We could interpolate using 'spatial statistics' (kriging)
- In fact we use a Kalman Filter instead
  - The Kalman Filter consists of two steps
  - A time update and a measurement update.

## Update Equations

---

- The time update equations project forward in time the current state and the error covariance estimate to give an *a priori* estimate for the next time step.
- The measurement update equations are responsible for the feedback, i.e. incorporating a new measurement into the *a priori* estimate.
- The time update equations are the predictor equations with the measurement update equations as the corrector equations.
- The final estimation algorithm resembles that of a predictor corrector algorithm for solving numerical problems



## Kalman Gain

---

- The data set was interpolated using a Kalman Filter
- The first task is to calculate the Kalman gain :

$$K = H^T R^{-1} (H^T R^{-1} H + S_a^{-1})^{-1}$$

- $S_a$  the priori covariance includes latitude, longitude and time components.

$$[S_a]_{ij} = a_0^2 a_1^{|\Delta_i - \Delta_j|} a_2^{|\Delta_1 - \Delta_2|}$$

- $R$  is the noise term (=  $\sigma^2 I$ ).
- $H$  is the measurement model (=1).

## Incorporating the Measurement

---

- The next step is to use the observations,  $y$ , and to generate an *a posteriori* state estimate by incorporating the measurement

$$\hat{x} = x_a + K(y - Hx_a)$$

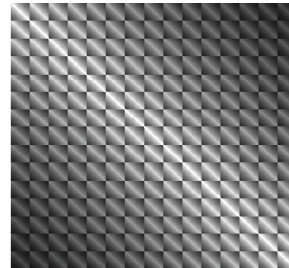
- $x_a = 0$  at first then  $x_{t-1}$  after.
- $y$  is the observed values.
- These equations operate on the linear version of the problem, and operates sequentially in time or space. In this case the interpolation is spatial.



## Step by step guide

- Initially we use the previous time step as an estimate for the missing data,  $x_a$
- As simplification, the measurement model,  $H$ , is taken as the identity matrix.
- A local Kalman filter is then used on a local spatial grid (15 by 15 degrees) as a method of 'spatial interpolation'. The weighting of the surrounding values uses the equation:

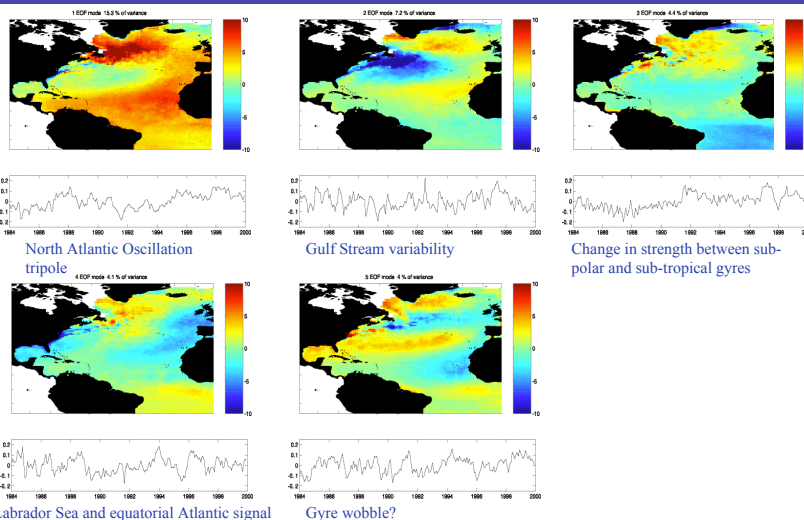
$$[S_a]_{ij} = a_0^2 a_1^{|\square_i - \square_j|} a_2^{|\square_1 - \square_2|}$$



$S_a$  the priori covariance for the 15 by 15° grid

- The generated value is used to 'replace' the missing value.

## North Atlantic modes of variability



## Probabilistic Principal Components

---

- Assume that the data we have is a linear combination of some 'hidden' or latent variables.

- Then

$$\mathbf{t} = \mathbf{W}\mathbf{x} + \boldsymbol{\mu} + \boldsymbol{\epsilon}$$

- Where  $\mathbf{t}$  is the observed data ( $n \times 1$ ),  $\mathbf{x}$  is the latent variables ( $p \times 1$ ),  $\mathbf{W}$  is a matrix ( $n \times p$ ),  $\boldsymbol{\mu}$  is the mean of  $\mathbf{t}$  and  $\boldsymbol{\epsilon}$  is an error term

## PPCA-2

---

- The  $\mathbf{x}$  are iid  $N(\mathbf{0}, \mathbf{I})$
- $\boldsymbol{\mu} \sim N(\mathbf{0}, \boldsymbol{\Sigma})$
- Then  $\mathbf{t} \sim N(\boldsymbol{\mu}, \mathbf{W}\mathbf{W}^t + \boldsymbol{\Sigma})$
- If  $\boldsymbol{\Sigma}$  is diagonal then this is Factor Analysis
- If we put  $\boldsymbol{\Sigma} = \boldsymbol{\Sigma}^2 \mathbf{I}$  then this becomes equivalent to Principal Components
- So PPCA enables us to calculate the first  $p$  principal components ( $(n-p)\boldsymbol{\Sigma}^2$  gives the variance not explained)

## The EM algorithm (Roweis)

---

- Expectation step

$$X = (C^T C)^{-1} C^T T$$

- Maximisation step

$$C^{new} = T X^T (X X^T)^{-1}$$

where:

$$C = W^T W$$

## The EM algorithm for Missing Data

---

- e-step  $X = (C^T C)^{-1} C^T Y$

$$Y = C X$$

- m-step  $C^{new} = Y X^T (X X^T)^{-1}$

- $Y$  is a matrix of the observed data ( $p \times n$ )
- $X$  is the matrix of the unknown states ( $k \times n$ )
- $C$  is the estimated rotated covariance matrix ( $p \times k$ )

These equations are iterated until they converge.

## Algorithm for PPCA with missing values

---

- Start with an initial guess of the covariance matrix.
- Cycle through a series of steps, replacing the missing values.
- Then re-estimate the covariance matrix.
- Continue until the values converge.

## How many EOF's to retain?

---

- 1) Look at the cumulative percentage variation explained.
- 2) Look at the size of the variance explained by each EOF - reject if less than one.
- 3) The scree plot.

However, none of these methods work for this technique.

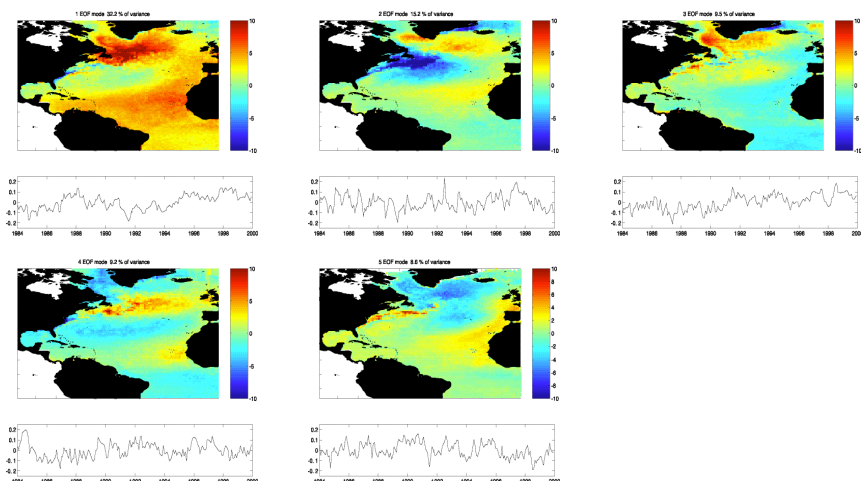
For this study we have decided to retain 9 EOF's

## Removing data to test the EM algorithm

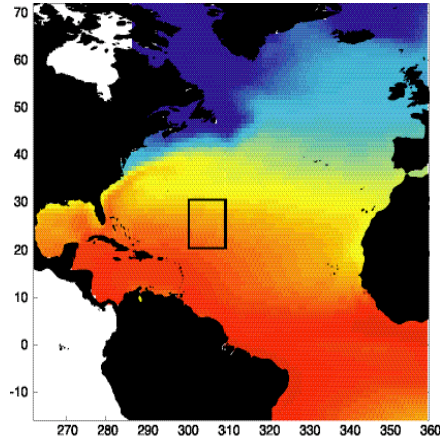
Two methods are used to test the EM algorithm

1. Data are removed randomly from the previously interpolated data set before the EM algorithm.
2. Data are removed from a block in the centre of the previously interpolated data set, randomly over time - to simulate large cloud bands.

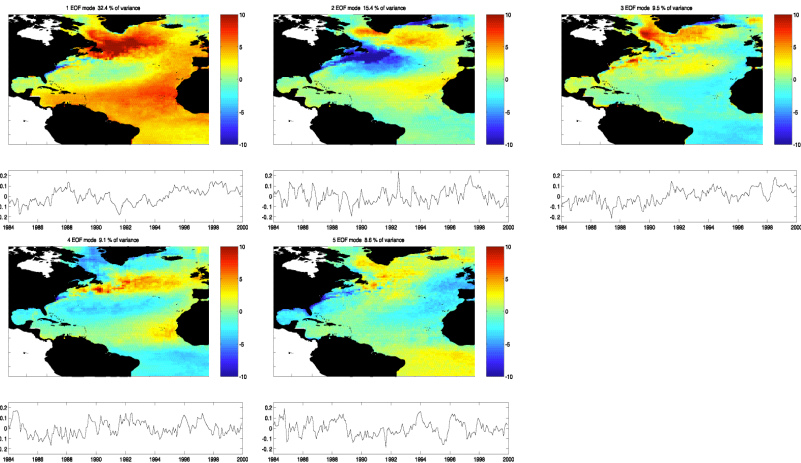
## First 5 EOF's using the EM Algorithm with 10% of data removed



# Removing data from a block



# First 5 EOF's using the EM Algorithm with 10% of data removed from a block



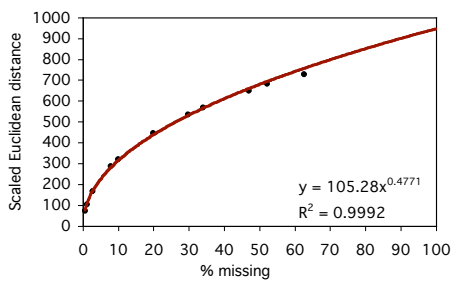
## EOF breakdown

The first five EOF's appeared unaffected by the missing values, until 30% of the data was removed

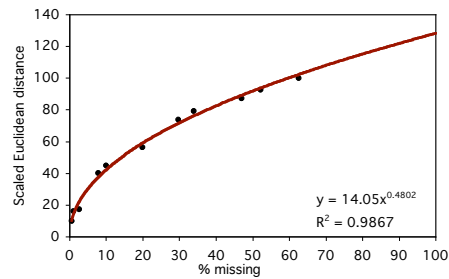
The reconstructed data was compared to the original data using the scaled Euclidean distance

$$D_s = \sqrt{\frac{\sum_{i=1}^N (x_i - y_i)^2}{N}}$$

## EOF Breakdown

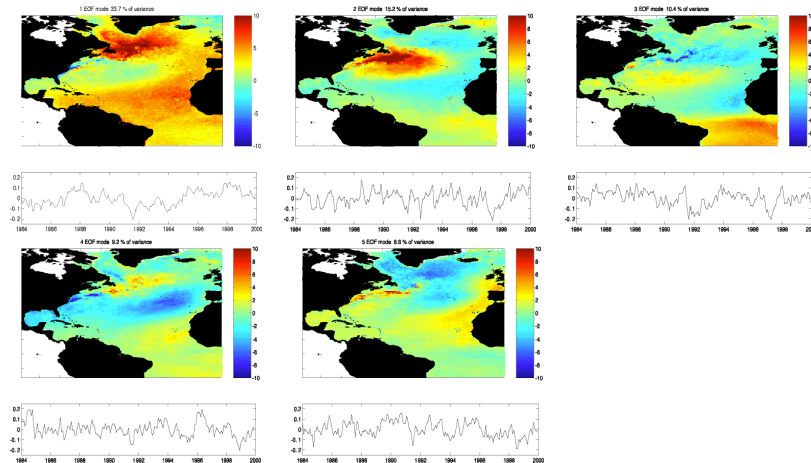


Random missing data



Block missing data

## First 5 EOF's using the EM Algorithm using uninterpolated data



## PPCA Summary

- The EM algorithm has produced similar results to those produced through standard procedures.
- Calculation of the EOF's was quicker, as the full covariance matrix was not calculated.
- The replacement of missing values was computationally more efficient than other interpolation methods.



## References

---

- Roweis, S., 1997: EM Algorithms for PCA and SPCA, *Neural Information Processing Systems 10 (NIPS'97)*, 626-632.
- Tipping, M.E. and C.M. Bishop, 1999: Probabilistic Principal component analysis. *Journal of the Royal Statistical Society, Series B*, 61(3), 611-622.
- Houseago-Stokes, R.E. and Challenor P.G. Using PPCA to estimate EOF's in the presence of missing values, Submitted to the *Journal of Atmospheric and Oceanic Technology*