

## Chapter 3 Linear Multivariate Statistical Analysis

As one often encounters datasets with more than a few variables, multivariate statistical techniques are needed to effectively extract the information contained in these datasets. In the environmental sciences, examples of multivariate datasets are ubiquitous --- the air temperatures recorded by all the weather stations around globe, the satellite infrared images composed of numerous small pixels, the gridded output from a general circulation model, ... The number of variables or time series from these datasets range from thousands to millions. Without a mastery of multivariate techniques, one is overwhelmed by these gigantic datasets. In this chapter, we review the principal component analysis method and its many variants, and the canonical correlation analysis method. These methods, using standard matrix techniques such as singular value decomposition, are relatively easy to use, but suffer from being linear, a limitation which will be lifted with neural network techniques in later chapters.

### 3.1 Principal component analysis (PCA)

#### 3.1.1 Geometric approach to PCA

We have a dataset with variables  $y_1, \dots, y_m$ . These variables have been sampled  $n$  times, e.g. the  $m$  variables could be  $m$  time series containing  $n$  observations in time. If  $m$  is a large number, we would like to capture the essence of  $y_1, \dots, y_m$  by a smaller set of variables  $z_1, \dots, z_k$  (i.e.  $k < m$ ; and hopefully  $k \ll m$  for truly large  $m$ ). This is the objective of principal component analysis (PCA), also called empirical orthogonal function (EOF) analysis in meteorology and oceanography. We first begin with a geometric approach, which is more intuitive than the standard eigenvector approach to PCA.

Let us start with only 2 variables  $y_1$  and  $y_2$ , as illustrated in Fig. 4. Clearly the bulk of the variance is along the axis  $z_1$ . If  $r_i$  is the distance between the  $i^{\text{th}}$  data point and the axis  $z_1$ , then the optimal  $z_1$  is found by minimizing  $\sum_{i=1}^n r_i^2$ .

Note that PCA treats all variables equally, whereas regression divides variables into independent and dependent variables, hence the straight line described by  $z_1$  is in general different from the regression line.

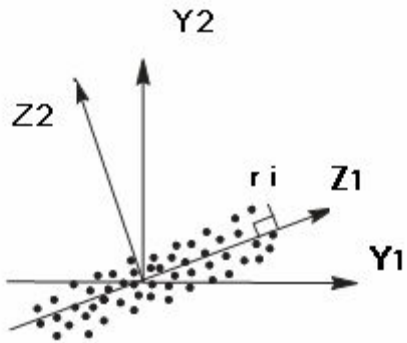


Fig. 4.1

### 3.1.2 Eigenvector approach to PCA

Taking the above example, a data point is transformed from its old coordinates  $(y_1, y_2)$  to new coordinates  $(z_1, z_2)$  via a rotation of the coordinate system (Fig. 4.2):

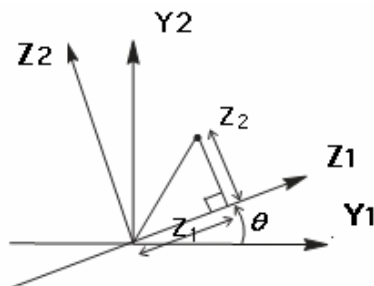
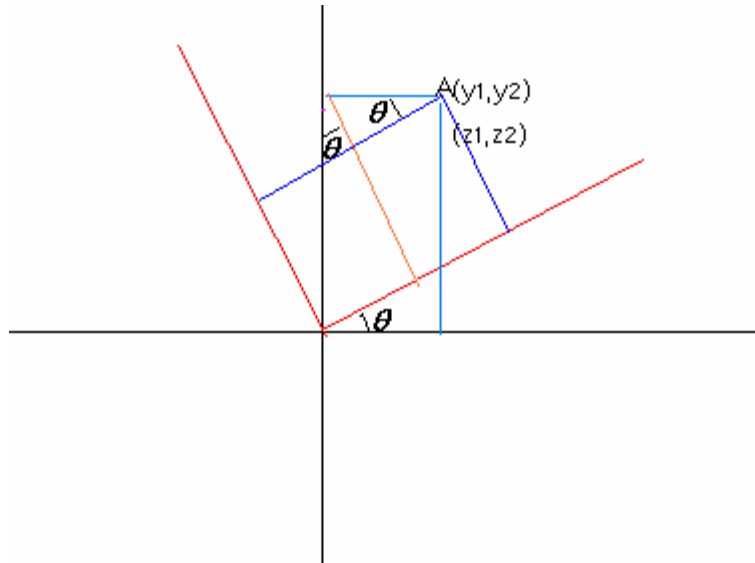


Fig. 4.2

$$\begin{aligned} z_1 &= y_1 \cos \theta + y_2 \sin \theta \\ z_2 &= -y_1 \sin \theta + y_2 \cos \theta \end{aligned} \quad (4.1)$$



In the general  $m$ -dimensional problem, we want to introduce new coordinates:

$$z_j = \sum_{l=1}^m e_{jl} y_l, j=1, \dots, m \quad (4.2)$$

The objective is to find

$$\mathbf{e}_1 = [e_{11}, \dots, e_{1m}]^T \quad (4.3)$$

which maximizes  $\text{var}(z_1)$ , i.e., find the coordinate transformation such that the variance of the dataset along the direction of the  $z_1$  axis is maximized.

$$\text{With } z_1 = \sum_{l=1}^m e_{1l} y_l = \mathbf{e}_1^T \mathbf{y}, \quad \mathbf{y} = [y_1 \dots y_m]^T,$$

We have

$$\text{var}(z_1) = \mathbf{E}[(z_1 - \bar{z}_1)(z_1 - \bar{z}_1)] = \mathbf{E}[\mathbf{e}_1^T (\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^T \mathbf{e}_1], \quad (4.4)$$

where we have used the vector property  $(a^T b)^T = b^T a$ . Thus,

$$\text{var}(z_1) = \mathbf{e}_1^T \mathbf{E}[(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^T] \mathbf{e}_1 = \mathbf{e}_1^T \mathbf{C} \mathbf{e}_1 \quad (4.5)$$

Where the covariance matrix  $\mathbf{C}$  is given by

$$\mathbf{C} = \mathbf{E}[(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^T] \quad (4.6)$$

Clearly, the larger is  $\|\mathbf{e}_1\|$ , the larger  $\text{var}(z_1)$  will be. Hence, we need to place a constraint on  $\|\mathbf{e}_1\|$  while we try to maximize  $\text{var}(z_1)$ . Let us impose a normalization constraint  $\|\mathbf{e}_1\|=1$ , i.e.

$$\mathbf{e}_1^T \mathbf{e}_1 = 1 \quad (4.7)$$

Thus our optimization problem is to find  $\mathbf{e}_1$  which maximize  $\mathbf{e}_1^T \mathbf{C} \mathbf{e}_1$ , subject to the constraint

$$\mathbf{e}_1^T \mathbf{e}_1 - 1 = 0 \quad (4.8)$$

The method of Lagrange multiplier s is commonly used to tackle optimization under constraints. Define the Lagrange function  $L$  by

$$L = \mathbf{e}_1^T \mathbf{C} \mathbf{e}_1 - \lambda(\mathbf{e}_1^T \mathbf{e}_1 - 1) \quad (4.9)$$

Where  $\lambda$  is a Lagrange multiplier.

To obtain  $\mathbf{e}_1$ , we ask for

$$\frac{\partial L}{\partial \mathbf{e}_1} = 0 \quad (4.10)$$

$$\mathbf{C} \mathbf{e}_1 - \lambda \mathbf{e}_1 = 0 \quad (4.11)$$

Which says that  $\lambda$  is an eigenvalue of the covariance matrix  $\mathbf{C}$ , which  $\mathbf{e}_1$  the eigenvector. Multiplying this equation by  $\mathbf{e}_1^T$  on the left, we obtain

$$\lambda = \mathbf{e}_1^T \mathbf{C} \mathbf{e}_1 = \text{var}(z_1) \quad (4.12)$$

Since  $\mathbf{e}_1^T \mathbf{C} \mathbf{e}_1$  is maximized, the so are  $\lambda$  and  $\text{var}(z_1)$ . The new coordinate  $z_1$ , called the principal component (PC), is found from (4.2).

Next, we want to find  $z_2$ --- our task is to find  $e_2$  which maximizes

$\text{var}(z_2) = e_2^T C e_2$ , subject to the constraint  $e_2^T e_2 = 1$ , and the constraint that  $z_2$  be uncorrelated with  $z_1$ , i.e., the covariance between  $z_1$  and  $z_2$  be zero,

$$\text{cov}(z_1, z_2) = 0. \quad (4.13)$$

As  $C = C^T$ , we can write

$$\begin{aligned} 0 &= \text{cov}(z_1, z_2) = \text{cov}(e_1^T y, e_2^T y) = E[e_1^T (y - \bar{y})(y - \bar{y}) e_2] \\ &= e_1^T C e_2 = e_2^T C e_1 = e_2^T \lambda_1 e_1 = \lambda_1 e_2^T e_1 = \lambda_1 e_1^T e_2 \end{aligned} \quad (4.14)$$

The orthogonal condition

$$e_2^T e_1 = 0 \quad (4.15)$$

can be used as a constraint in place of (4.13).

Upon introducing another Lagrange multiplier  $\gamma$ , we want to find an  $e_2$  which gives a stationary point of the Lagrange function  $L$ ,

$$L = e_2^T C e_2 - \lambda(e_2^T e_2 - 1) - \gamma e_2^T e_1 \quad (4.16)$$

$$\frac{\partial L}{\partial e_1} = 0, \quad \frac{\partial L}{\partial e_2} = 0$$

$$C e_2 - \lambda e_2 = 0 \quad (4.17)$$

Once again  $\lambda$  is an eigenvalue of the covariance matrix  $C$ , which  $e_2$  the eigenvector. As

$$\lambda = e_2^T C e_2 = \text{var}(z_2) \quad (4.18)$$

which is maximized, this  $\lambda = \lambda_2$  is as large as possible with  $\lambda_2 < \lambda_1$ . (The case  $\lambda_2 = \lambda_1$  is degenerate and will be discussed later). Hence,  $\lambda_2$  is the second largest eigenvalue of  $C$ , with  $\lambda_2 = \text{var}(z_2)$ . This process can be repeated for  $z_3, z_4, \dots$

So far,  $\mathbf{C}$  is the data covariance matrix, but it can also be the data correlation matrix, if one prefers correlation over covariance. In combined PCA, where two or more variables with different units are combined into one large data matrix for PCA --- e.g. finding the PCA modes of the combined sea surface temperature data and the sea level pressure data --- then one needs to normalize the variables so  $\mathbf{C}$  is the correlation matrix.

So, the general procedure of PCA analysis for dataset  $\mathbf{y}$  is as below:

- (1) calculating covariance matrix (or correlation matrix)  $\mathbf{C} = \{C_{ij}\}$

$$C_{i,j} = \sum_{k=1}^N (x_{ik} - \bar{x}_i)(x_{jk} - \bar{x}_j) \quad (i=1, \dots, M, j=1, \dots, M)$$

where  $M$  denote the number of variables (or grids), and  $N$ , the length of samples.

- (2) calculating the eigenvalues and eigenvectors of  $\mathbf{C}$

$$\mathbf{C}\mathbf{e} - \lambda\mathbf{e} = 0$$

where  $\mathbf{e} = \{\mathbf{e}_1, \dots, \mathbf{e}_L\}$ ,  $\lambda = \{\lambda_1, \dots, \lambda_L\}$   $L = \min\{M, N\}$

- (3) Usually  $\lambda_1 > \lambda_2 > \dots > \lambda_L$  but sometimes, the outputs are in reverse order. In this case, the first eigenvector corresponding with  $\lambda_L$ .

- (3) calculating the PCs, i.e.,  $\mathbf{z} = \mathbf{e}^T \mathbf{y}$

### 3.1.3 Real and complex data

In general, for  $\mathbf{y}$  real,

$$\mathbf{C} = \mathbf{E}[(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^T] \quad (4.19)$$

implies that  $\mathbf{C}^T = \mathbf{C}$ , i.e.,  $\mathbf{C}$  is a real, symmetric matrix. A positive semi-defined matrix  $\mathbf{A}$  is defined by the property that for any  $\mathbf{v} \neq 0$ , it follows that  $\mathbf{v}^T \mathbf{A} \mathbf{v} \geq 0$ . From the definition of  $\mathbf{C}$  (4.5), it is clear that  $\mathbf{v}^T \mathbf{C} \mathbf{v} \geq 0$  is satisfied. Hence  $\mathbf{C}$  is a real, symmetric, positive semi-definite matrix.

If  $\mathbf{y}$  is complex, then

$$\mathbf{C} = \mathbf{E}[(\mathbf{y} - \bar{\mathbf{y}})(\mathbf{y} - \bar{\mathbf{y}})^{T*}] \quad (4.20)$$

With complex conjugation denoted by the superscript asterisk. As  $\mathbf{C}^{T*} = \mathbf{C}$ ,  $\mathbf{C}$  is a Hermitian matrix. It is also a positive semi-definite matrix.

Theorems on Hermitian matrix, positive semi-definite matrices tell us that:  $\mathbf{C}$  has real eigenvalues:

$$\lambda_1 > \lambda_2 > \dots > \lambda_L \geq 0, \quad \sum_{j=1}^L \lambda_j = \text{var}(\mathbf{y}) \quad (4.21)$$

### 3.1.4 Orthogonality relations

Thus PCA amounts to finding the eigenvectors and eigenvalues of  $\mathbf{C}$ . The orthogonal eigenvectors then provide a basis, i.e., the data  $\mathbf{y}$  can be expanded in terms of the eigenvectors  $\mathbf{e}_j$ :

$$\mathbf{y} - \bar{\mathbf{y}} = \sum_{j=1}^m \mathbf{a}_j(t) \mathbf{e}_j \quad (4.22)$$

where  $\mathbf{a}_j(t)$  are the expansion coefficients. To obtain  $\mathbf{a}_j(t)$ , left multiply the above equation by  $\mathbf{e}_i^T$ , and use the orthogonal relation of the eigenvectors,

$$\mathbf{e}_i^T \mathbf{e}_j = \delta_{ij} \quad (4.23)$$

to get

$$\mathbf{a}_j(t) = \mathbf{e}_j^T (\mathbf{y} - \bar{\mathbf{y}}), \quad (4.24)$$

i.e.,  $\mathbf{a}_j(t)$  is obtained by the projection of the data vector  $\mathbf{y} - \bar{\mathbf{y}}$  onto the eigenvector  $\mathbf{e}_j$ , as the right hand side of this equation is simply a dot product between the two vectors.  $\mathbf{a}_j(t)$  are usually called PCs, and  $\mathbf{e}_j$  eigenvectors or EOFs.

There are two important properties:

(1) The expansion  $\sum_{j=1}^k a_j(t) e_j(x)$  explains more of the variances of the data than any other linear combination  $\sum_{j=1}^k b_j(t) f_j(x)$ . Thus PCA provides the most efficient way to compress data.

(2) The time series in the set  $\{a_j\}$  are uncorrelated. We can write

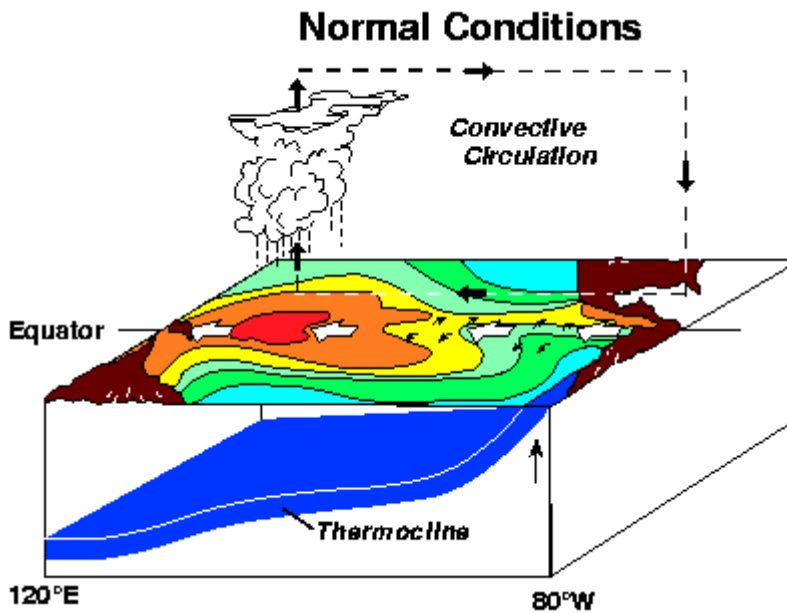
$$a_j(t) = e_j^T (y - \bar{y}) = (y - \bar{y}) e_j$$

For  $i \neq j$ ,

$$\begin{aligned} \text{cov}(a_i, a_j) &= E[e_i^T (y - \bar{y})(y - \bar{y}) e_j] = e_i^T E[(y - \bar{y})(y - \bar{y})^T] e_j \\ &= e_i^T C e_j = e_i^T \lambda_j e_j = \lambda_j e_i^T e_j = 0 \end{aligned} \quad (4.25)$$

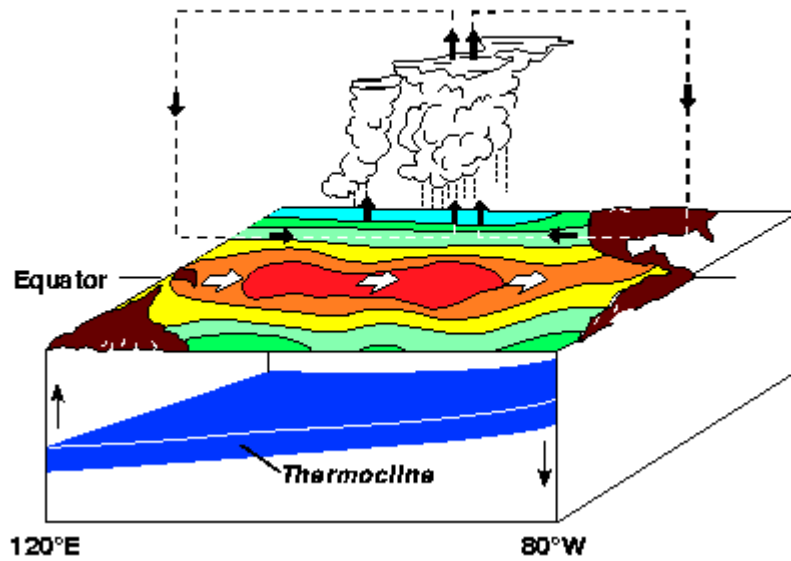
Hence PCA extracts the uncorrelated modes of variability of the data field. Note that no correlation between  $(a_i, a_j)$  only means no linear relation between the two, there may still be nonlinear relation between them, which can be extracted by the nonlinear PCA method using neural network.

### 3.1.5 An Example: PCA of the tropical Pacific climate variability

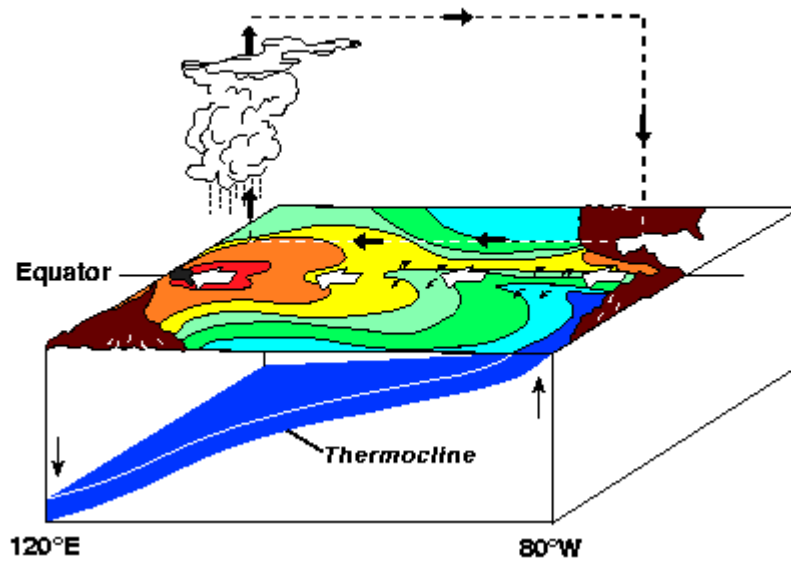


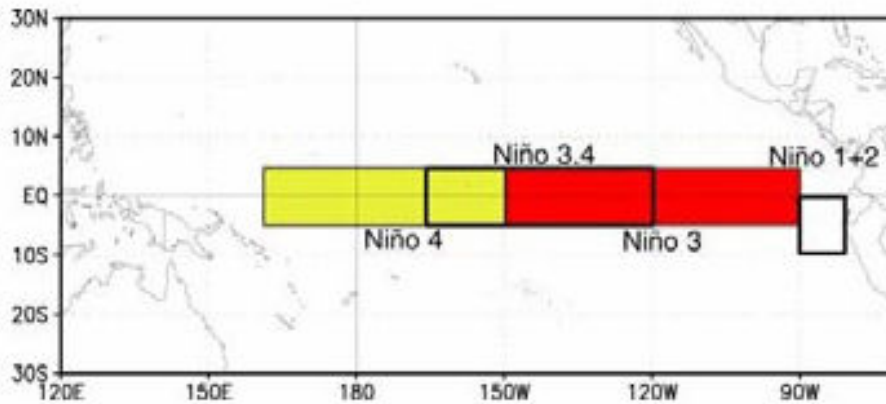


### El Niño Conditions



### La Niña Conditions





Let us study the monthly tropical Pacific SST from NOAA (National Oceanic and atmospheric administration) for the period January, 1950 to August 2000. The SST field has 2 spatial dimensions, but can easily be rearranged into the form of  $y(t)$  for PCA analysis. Fig.4.3 is the spatial pattern for the first 3 modes (accounting for 51.8%, 10.1% and 7.3% respectively, of the total SST variance). Fig.4.4 are PCs corresponding with the first 3 modes.

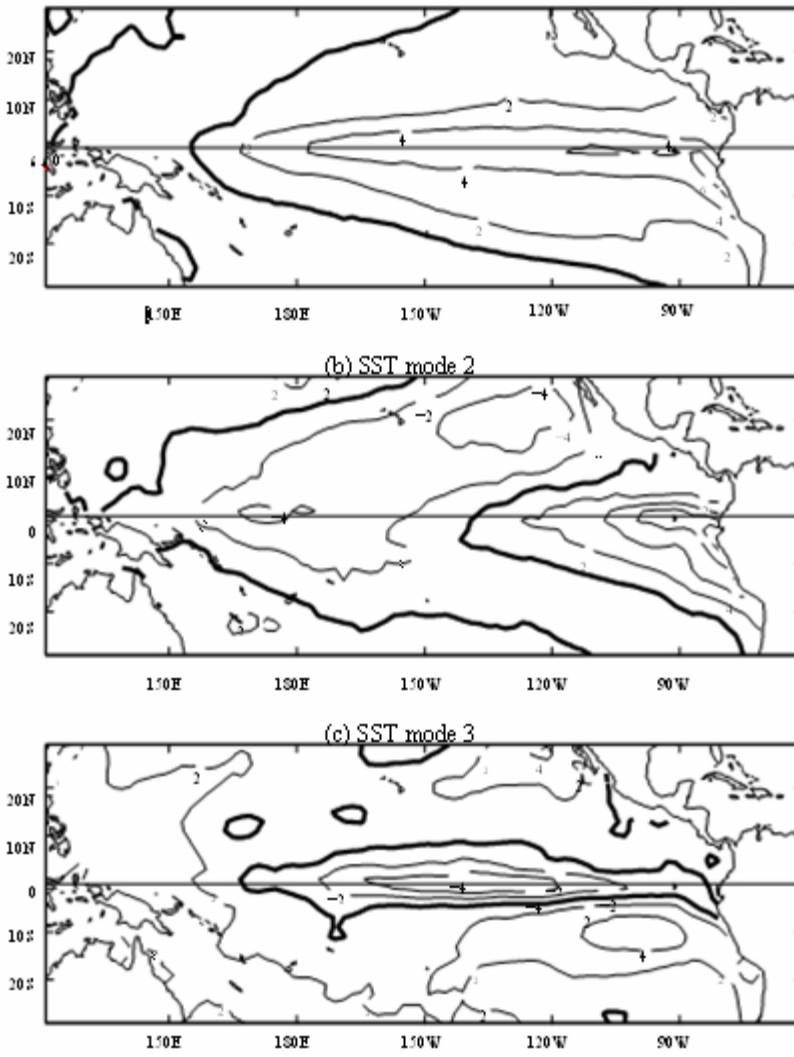


Fig4.3

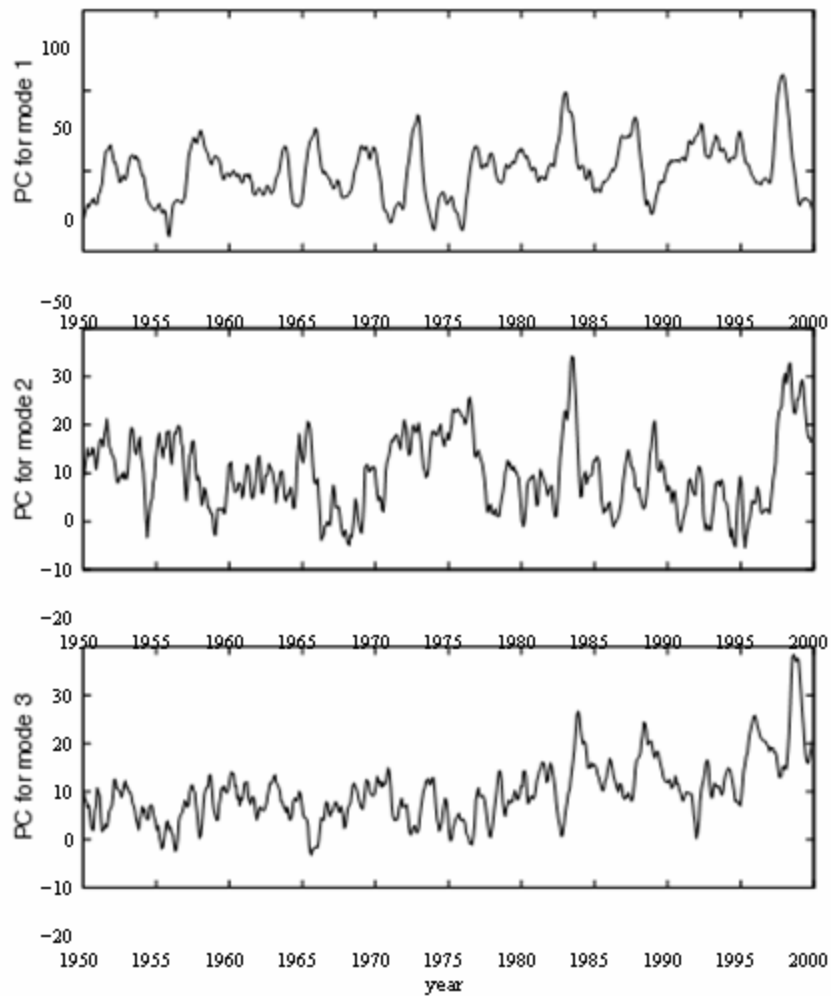


Fig.4.4

Mode1 (Fig. 4.3a) shows the largest SST anomalies occurring in the eastern and central equatorial Pacific. The PC1 (Fig4.4) can be used as an index for El Nino/La Nina.

Mode2 (Fig. 4.3b) has, along the equator, positive anomalies near the east and negative anomalies further west.

Mode3 (Fig. 4.3c) shows the largest anomaly occurring in the central equatorial Pacific, and the PC shows a rising trend after the mid 1970s.

### 3.1.6 Scaling the PCs and eigenvectors

There are various options for the scaling of the PCs  $\{\mathbf{a}_j(t)\}$  and the eigenvectors  $\{\mathbf{e}_j\}$ . One can introduce an arbitrary scale factor  $\alpha$ ,

$$\mathbf{a}_j' = \frac{1}{\alpha} \mathbf{a}_j \quad \mathbf{e}_j' = \alpha \mathbf{e}_j \quad (4.26)$$

so that

$$\mathbf{y} - \bar{\mathbf{y}} = \sum_j \mathbf{a}_j' \mathbf{e}_j' \quad (4.27)$$

Our choice for the scaling has so far been

$$\mathbf{e}_i^T \mathbf{e}_j = \delta_{ij} \quad (4.28)$$

which was the choice of Lorenz. The variance of the original data  $\mathbf{y}$  is then contained in  $\{\mathbf{a}_j(t)\}$ , with

$$\text{var}(\mathbf{y}) = \mathbf{E} \left[ \sum_{j=1}^m \mathbf{a}_j^2 \right] \quad (4.29)$$

Another common choice is Hotelling's original choice

$$\mathbf{a}_j' = \frac{1}{\sqrt{\lambda_j}} \mathbf{a}_j, \quad \mathbf{e}_j' = \sqrt{\lambda_j} \mathbf{e}_j \quad (4.30)$$

whence

$$\text{var}(\mathbf{y}) = \sum_{j=1}^m \lambda_j = \sum_{j=1}^m \|\mathbf{e}_j'\|^2$$

$$\text{From (4.25), } \text{cov}(\mathbf{a}_i', \mathbf{a}_j') = \frac{1}{\sqrt{\lambda_i} \sqrt{\lambda_j}} \text{cov}(\mathbf{a}_i, \mathbf{a}_j) = \frac{\lambda_j}{\sqrt{\lambda_i} \sqrt{\lambda_j}} \mathbf{e}_i^T \mathbf{e}_j = \delta_{ij} \quad (4.31)$$

The variance of the original data is now contained in  $\{\mathbf{e}_j(t)\}$  instead. In sum, regardless of the arbitrary scale factor, the PCA eigenvectors are orthogonal and the PCs are uncorrelated.

If  $\tilde{y}_i$  is  $y_i$  with mean removed and normalized by standard deviation, then one can show that the correlation

$$\rho(\mathbf{a}_j'(t), \tilde{\mathbf{y}}_l) = e_{jl}' \quad (4.32)$$

the  $l$ th element of  $e_j'$ . Hence the  $l$ th element of  $e_j'$  conveniently provides the correlation between the PC  $\mathbf{a}_j'$  and the standardized variable  $\tilde{y}_l$ .

### 3.1.7 Degeneracy of Eigenvalues

A degenerate case arises when  $\lambda_i = \lambda_j$ . When two eigenvalues are equal, the eigenvectors are not unique.

A simple example of degeneracy is illustrated by a propagating plane wave,

$$h(x, y, t) = A \cos(ky - \omega t) \quad (4.33)$$

Which can be expressed in terms of two standing waves:

$$h(x, y, t) = A \cos(ky) \cos(\omega t) + A \sin(ky) \sin(\omega t) \quad (4.34)$$

If we perform PCA on  $h(x, y, t)$ , we get two modes with equal eigenvalues. To see this, note that in the  $x$ - $y$  plane,  $\cos(ky)$  and  $\sin(ky)$  are orthogonal, while  $\cos(\omega t)$  and  $\sin(\omega t)$  are uncorrelated, so (4.34) satisfies the properties of PCA modes in that the eigenvectors are orthogonal and the PCs are uncorrelated. As (4.34) is a PCA decomposition, with the two modes both having the same amplitude  $A$ , hence the eigenvalues  $\lambda_1 = \lambda_2$ , and the case is degenerate. Thus propagating waves in the data leads to degeneracy in the eigenvalues. If one finds eigenvalues of very similar magnitudes from a PCA analysis, that implies near degeneracy and there may be propagating waves in the data.

### 3.1.8 A smaller covariance matrix

Let the data matrix be

$$\mathbf{Y} = \begin{bmatrix} y_{11} & \cdots & y_{1n} \\ \cdots & \cdots & \cdots \\ y_{m1} & \cdots & y_{mn} \end{bmatrix} \quad (4.35)$$

Where  $m$  is the number of spatial points and  $n$  the number of time points.

Assuming the temporal mean has been removed, then covariance matrix

$$\mathbf{C} = \frac{1}{n} \mathbf{Y} \mathbf{Y}^T \quad \text{and} \quad \mathbf{C}' = \frac{1}{n} \mathbf{Y}^T \mathbf{Y} \quad (4.36)$$

$\mathbf{C}$  is a  $m \times m$  matrix, and  $\mathbf{C}'$  is a  $n \times n$  matrix.

In most problems, the size of the two matrices are very different. For instance, for global  $5^\circ \times 5^\circ$  monthly sea level pressure data collected over 50 years, the total number of spatial grid points is  $m=2592$  while the number of time points is  $n=600$ . Obviously, it will be much cheaper to solve the eigen problem for  $\mathbf{C}'$  than for  $\mathbf{C}$ .

The matrix theory says:  $\mathbf{C}$  and  $\mathbf{C}'$  have same eigenvalues. The question is now how to get eigenvectors of  $\mathbf{C}$  from eigenvectors of  $\mathbf{C}'$ ?

$$\mathbf{C}' = \frac{1}{n} \mathbf{Y}^T \mathbf{Y}$$

$$\frac{1}{n} \mathbf{Y}^T \mathbf{Y} \mathbf{v}_j = \lambda_j \mathbf{v}_j \quad (4.37)$$

where  $\mathbf{v}_j$  and  $\lambda_j$  are eigenvectors and eigenvalues of  $\mathbf{C}'$ .

Multiplying  $\mathbf{Y}$  on both sides of (4.36), we have

$$\mathbf{Y} \frac{1}{n} \mathbf{Y}^T \mathbf{Y} \mathbf{v}_j = \mathbf{Y} \lambda_j \mathbf{v}_j \quad (4.38)$$

$$\left( \frac{1}{n} \mathbf{Y} \mathbf{Y}^T \right) (\mathbf{Y} \mathbf{v}_j) = \lambda_j (\mathbf{Y} \mathbf{v}_j) \quad (4.39)$$

Denoting

$$\mathbf{e}_j = \mathbf{Y} \mathbf{v}_j, \quad (4.40)$$

we have

$$\mathbf{C} \mathbf{e}_j = \lambda_j \mathbf{e}_j \quad (4.41)$$

(4.41) is just the eigen equation for  $\mathbf{C}$ , meaning  $\mathbf{e}_j$  is an eigenvector for  $\mathbf{C}$ .

In summary, solving the eigen problem for the smaller matrix  $\mathbf{C}'$  yields the eigenvalues  $\{\lambda_j\}$  and eigenvectors  $\{\mathbf{v}_j\}$ . The eigenvectors  $\{\mathbf{e}_j\}$  for the bigger matrix  $\mathbf{C}$  are then obtained from (4.40)  $\mathbf{e}_j = \mathbf{Y}\mathbf{v}_j$ .

### 3.1.9 Singular value decomposition

Instead of solving the eigen problem of the data covariance matrix  $\mathbf{C}$ , a computationally more efficient way to perform PCA is via Singular Value Decomposition (SVD) of the  $m \times n$  data matrix  $\mathbf{Y}$  given by (4.35). Without loss of generality, we can assume  $m \geq n$ , then the SVD Theorem says that

$$\mathbf{Y} = \mathbf{E}\mathbf{S}\mathbf{F}^T \quad (4.42)$$

$$\mathbf{E} : m \times m; \quad \mathbf{S} : m \times n; \quad \mathbf{F} : n \times n$$

$\mathbf{E}$  and  $\mathbf{F}$  are orthonormal matrices, i.e., they satisfy

$$\mathbf{E}^T \mathbf{E} = \mathbf{I}, \quad \mathbf{F}^T \mathbf{F} = \mathbf{I}, \quad (4.43)$$

Where  $\mathbf{I}$  is the identity matrix. The leftmost  $n$  columns of  $\mathbf{E}$  contain the  $n$  left singular vectors, and then columns of  $\mathbf{F}$  the  $n$  right singular vectors, while the diagonal elements of  $\mathbf{S}$  are the singular values.

The covariance matrix  $\mathbf{C}$  can be rewritten as

$$\mathbf{C} = \frac{1}{n} \mathbf{Y}\mathbf{Y}^T = \frac{1}{n} \mathbf{E}\mathbf{S}\mathbf{S}^T \mathbf{E}^T \quad (4.44)$$

$$\text{The matrix } \frac{1}{n} \mathbf{S}\mathbf{S}^T = \mathbf{\Lambda} \quad (4.45)$$

is diagonal and zero everywhere, except in the upper left  $n \times n$  corner, containing  $\frac{1}{n} s_{(i,i)}^2$  ---  $s_{(i,i)}$  is singular values

Right multiply Eq. (4.44) by  $\mathbf{E}$  gives



$$CE = E\Lambda \quad (4.46)$$

(4.46) is a standard eigen equation, with  $E$  is the eigenvectors and  $\Lambda$  is the eigenvalues. So, we can use SVD to derive eigenvectors and eigenvalues with the relation:  $\lambda_j = \frac{1}{n} s(j, j)^2$ .

SVD approach to PCA is at least twice as fast as the eigen approach. So, SVD is in particular useful for large datasets. Matlab program for SVD is [svd](#).

### 3.1.20 Significance tests

In practice, the higher PCA modes, which basically contain noise, are rejected. How does one decide how many modes to retain? There are some “rules of thumb”. One of the simplest approach is to plot the eigenvalues  $\lambda_j$  as a function of the mode number  $j$ . Hopefully, from the plot, one finds an abrupt transition from large eigenvalues to small eigenvalues around mode number  $m$ . One can then retain the first  $m$  modes. Alternatively, the Kaiser test rejects the modes with eigenvalues  $\lambda$  less than the mean value  $\bar{\lambda}$ .

Computationally more involved is the Monte Carlo test, which involves setting up random data matrices  $R_k$  ( $k = 1, \dots, K$ ), of the same size as the data matrix  $Y$ . The random elements are normally distributed, with the variance of the random data matching the variance of the actual data. PCA is performed on each of the random matrices, yielding eigenvalues  $\lambda^{(k)}_j$ . Assume for each  $k$ , the set of eigenvalues are sorted in descending order. For each  $j$ , one examines the distribution of the  $K$  values of  $\lambda^{(k)}_j$ , and finds the level  $\lambda_{0.05}$ , which is exceeded only by 5% of the  $\lambda^{(k)}_j$  values. The eigenvalues  $\lambda_j$  from  $Y$  which failed to rise above this  $\lambda_{0.05}$  level are then rejected.

Since the Monte Carlo method performs PCA on  $K$  matrices and  $K$  is typically about 100, it can be costly for large data matrices.

## 3.2 Canonical correlation analysis (CCA)

Given a set of variables  $\{y_j\}$ , PCA finds the linear modes accounting for the maximum amount of variance in the dataset. When there are two sets of variables  $\{x_i\}$  and  $\{y_j\}$ , *canonical correlation analysis* (CCA), first introduced by \*\*\*\*Hotelling (1936), finds the modes of maximum correlation between  $\{x_i\}$  and  $\{y_j\}$ , rendering CCA a standard tool for discovering linear relations between two fields. CCA is a generalization of the Pearson correlation between two variables  $x$  and  $y$  to two sets of variables  $\{x_i\}$  and  $\{y_j\}$ . Thus CCA can be viewed as a “doubled-barreled PCA”. A variant of the CCA method finds the modes of maximum *covariance* between  $\{x_i\}$  and  $\{y_j\}$ — this variant is called the *maximum covariance analysis* (MCA) by von Storch and Zwiers (1999), and because it uses the SVD matrix technique, it is also simply called the SVD (singular value decomposition) method by other researchers, though this name is confusing as it is used to denote both a matrix technique and a multivariate statistical technique.

In PCA, one finds a linear combination of the  $y_j$  variables, i.e.  $\mathbf{e}_1^T \mathbf{y}$ , which has the largest variance (subject to  $\|\mathbf{e}_1\| = 1$ ). Next, one finds  $\mathbf{e}_2^T \mathbf{y}$  with the largest variance, but with  $\mathbf{e}_2^T \mathbf{y}$  uncorrelated with  $\mathbf{e}_1^T \mathbf{y}$ ; and similarly for the higher modes.

In CCA, one finds  $\mathbf{f}_1$  and  $\mathbf{g}_1$ , so that the correlation between  $\mathbf{f}_1^T \mathbf{x}$  and  $\mathbf{g}_1^T \mathbf{y}$  is maximized. Next find  $\mathbf{f}_2$  and  $\mathbf{g}_2$  so that the correlation between  $\mathbf{f}_2^T \mathbf{x}$  and  $\mathbf{g}_2^T \mathbf{y}$  is maximized, with  $\mathbf{f}_2^T \mathbf{x}$  and  $\mathbf{g}_2^T \mathbf{y}$  uncorrelated with both  $\mathbf{f}_1^T \mathbf{x}$  and  $\mathbf{g}_1^T \mathbf{y}$ . And so forth for the higher modes.

### 3.2.1 CCA theory

Consider two datasets

$$\mathbf{x}(t) = x_{il}, \quad i = 1, \dots, n_x, \quad l = 1, \dots, n_t, \quad (4.61)$$

and

$$\mathbf{y}(t) = y_{jl}, \quad j = 1, \dots, n_y, \quad l = 1, \dots, n_t. \quad (4.62)$$

i.e.,  $\mathbf{x}$  and  $\mathbf{y}$  need not have the same spatial dimensions, but need the same time dimension  $n_t$ . Assume  $\mathbf{x}$  and  $\mathbf{y}$  have zero means. Let

$$\mathbf{u} = \mathbf{f}^T \mathbf{x}, \quad \mathbf{v} = \mathbf{g}^T \mathbf{y}. \quad (4.63)$$

The correlation

$$\rho = \frac{\text{cov}(\mathbf{u}, \mathbf{v})}{\sqrt{\text{var}(\mathbf{u}) \text{var}(\mathbf{v})}} = \frac{\text{cov}(\mathbf{f}^T \mathbf{x}, \mathbf{g}^T \mathbf{y})}{\sqrt{\text{var}(\mathbf{f}^T \mathbf{x}) \text{var}(\mathbf{g}^T \mathbf{y})}} = \frac{\mathbf{f}^T \text{cov}(\mathbf{x}, \mathbf{y}) \mathbf{g}}{\sqrt{\text{var}(\mathbf{f}^T \mathbf{x}) \text{var}(\mathbf{g}^T \mathbf{y})}}, \quad (4.64)$$

Where we have invoked

$$\text{cov}(\mathbf{f}^T \mathbf{x}, \mathbf{g}^T \mathbf{y}) = \mathbb{E}(\mathbf{f}^T \mathbf{x} (\mathbf{g}^T \mathbf{y})^T) = \mathbb{E}(\mathbf{f}^T \mathbf{x} \mathbf{y}^T \mathbf{g}) = \mathbf{f}^T \mathbb{E}[\mathbf{x} \mathbf{y}^T] \mathbf{g}. \quad (4.65)$$

We want  $u$  and  $v$ , the two canonical variates or canonical correlation coordinates, to have maximum correlation between them, i.e.,  $f$  and  $g$  are chosen to maximize  $\rho$ . We are of course free to normalize  $f$  and  $g$  as we like, because if  $f$  and  $g$  maximize  $\rho$ , so will  $\alpha f$  and  $\beta g$ , for any nonzero  $\alpha$  and  $\beta$ . We choose the normalization condition:

$$\text{var}(\mathbf{f}^T \mathbf{x}) = 1 = \text{var}(\mathbf{g}^T \mathbf{y}). \quad (4.66)$$

Since

$$\text{var}(\mathbf{f}^T \mathbf{x}) = \text{cov}(\mathbf{f}^T \mathbf{x}, \mathbf{f}^T \mathbf{x}) = \mathbf{f}^T \text{cov}(\mathbf{x}, \mathbf{x}) \mathbf{f} \equiv \mathbf{f}^T \mathbf{C}_{xx} \mathbf{f} \quad (4.67)$$

and

$$\text{var}(\mathbf{g}^T \mathbf{y}) = \mathbf{g}^T \mathbf{C}_{yy} \mathbf{g} \quad (4.68)$$

So, we have

$$\mathbf{f}^T \mathbf{C}_{xx} \mathbf{f} = 1, \quad \mathbf{g}^T \mathbf{C}_{yy} \mathbf{g} = 1. \quad (4.69)$$

With (4.66), (4.64) reduces to

$$\rho = \mathbf{f}^T \mathbf{C}_{xy} \mathbf{g}, \quad (4.70)$$

Where  $\mathbf{C}_{xy} = \text{cov}(\mathbf{x}, \mathbf{y})$

The problem is to maximize (4.70) subject to constraints (4.69). We will again use the method of Lagrange multipliers, where we impose the constraints into the Lagrange function  $L$ ,

$$L = \mathbf{f}^T \mathbf{C}_{xy} \mathbf{g} + \alpha (\mathbf{f}^T \mathbf{C}_{xx} \mathbf{f} - 1) + \beta (\mathbf{g}^T \mathbf{C}_{yy} \mathbf{g} - 1), \quad (4.71)$$

where  $\alpha$  and  $\beta$  are the unknown Lagrange multipliers. To find the stationary points of  $L$ , we need

$$\frac{\partial L}{\partial \mathbf{f}} = \mathbf{C}_{xy} \mathbf{g} + 2\alpha \mathbf{C}_{xx} \mathbf{f} = 0, \quad (4.72)$$

$$\frac{\partial L}{\partial \mathbf{g}} = \mathbf{C}_{xy}^T \mathbf{f} + 2\beta \mathbf{C}_{yy} \mathbf{g} = 0.$$

Hence

$$\mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \mathbf{g} = -2\alpha \mathbf{f}, \quad (4.73)$$

$$\mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^T \mathbf{f} = -2\beta \mathbf{g}. \quad (4.74)$$

Substituting (4.74) into (4.73) yields

$$\mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^T \mathbf{f} \equiv \mathbf{M}_f \mathbf{f} = \lambda \mathbf{f}, \quad (4.75)$$

With  $\lambda = 4\alpha\beta$ . Similarly, substituting (4.73) into (4.74) yields

$$\mathbf{C}_{yy}^{-1} \mathbf{C}_{xy}^T \mathbf{C}_{xx}^{-1} \mathbf{C}_{xy} \mathbf{g} \equiv \mathbf{M}_g \mathbf{g} = \lambda \mathbf{g}. \quad (4.76)$$

Both these equations can be viewed as eigenvalue equations, with  $\mathbf{M}_f$  and  $\mathbf{M}_g$  sharing the same non-zero eigenvalues  $\lambda$ . As  $\mathbf{M}_f$  and  $\mathbf{M}_g$  are known from the data,  $\mathbf{f}$  can be found by solving the eigenvalue problem (4.75).  $\beta \mathbf{g}$  can then be obtained from (4.74). Since  $\beta$  is unknown, the magnitude of  $\mathbf{g}$  is unknown, and the normalization conditions (4.69) are used to determine the magnitude of  $\mathbf{g}$ . Alternatively, one can use (4.76) to solve for  $\mathbf{g}$  first, then obtain  $\mathbf{f}$  from (4.73), and the normalization condition (4.69). The matrix  $\mathbf{M}_f$  is of dimension  $n_x \times n_x$ , while  $\mathbf{M}_g$  is  $n_y \times n_y$ , so one usually picks the smaller of the two to solve the eigenvalue problem.

From (4.70),

$$\rho^2 = \mathbf{f}^T \mathbf{C}_{xy} \mathbf{g} \mathbf{g}^T \mathbf{C}_{xy}^T \mathbf{f} = 4\alpha\beta (\mathbf{f}^T \mathbf{C}_{xx} \mathbf{f}) (\mathbf{g}^T \mathbf{C}_{yy} \mathbf{g}), \quad (4.77)$$

Where (4.72) has been invoked. From (4.69), (4.77) reduces to

$$\rho^2 = \lambda. \quad (4.78)$$

The eigenvalue problems (4.75) and (4.76) yield  $n$  number of  $\lambda_s$ , with  $n = \min(n_x, n_y)$ . Assuming the  $\lambda_s$  to be all distinct and nonzero, we have for each  $\lambda_j$  ( $j=1, \dots, n$ ), a pair of eigenvectors,  $\mathbf{f}_j$  and  $\mathbf{g}_j$ , and a pair of canonical variates,  $u_j$  and  $v_j$ , with correlation  $\rho_j = \sqrt{\lambda_j}$  between the two.

Let us write the forward mapping from the variables  $\mathbf{x}(t)$  and  $\mathbf{y}(t)$  to the canonical variates  $\mathbf{u}(t) = [u_1(t), \dots, u_n(t)]^T$  and  $\mathbf{v}(t) = [v_1(t), \dots, v_n(t)]^T$  as

$$\mathbf{u} = [\mathbf{f}_1^T \mathbf{x}, \dots, \mathbf{f}_n^T \mathbf{x}]^T = \mathcal{F}^T \mathbf{x}, \quad \mathbf{v} = \mathcal{G}^T \mathbf{y} \quad (4.79)$$

Next, we need to find the inverse mapping from  $\mathbf{u}(t) = [u_1(t), \dots, u_n(t)]^T$  and  $\mathbf{v}(t) = [v_1(t), \dots, v_n(t)]^T$  to the original variables  $\mathbf{x}$  and  $\mathbf{y}$ . Let

$$\mathbf{x} = \mathbf{F} \mathbf{u}, \quad \mathbf{y} = \mathbf{G} \mathbf{v}. \quad (4.80)$$

We note that

$$\text{cov}(\mathbf{x}, \mathbf{u}) = \text{cov}(\mathbf{x}, \mathcal{F}^T \mathbf{x}) = \text{E}[\mathbf{x}(\mathcal{F}^T \mathbf{x})^T] = \text{E}[\mathbf{x} \mathbf{x}^T \mathcal{F}] = \mathbf{C}_{xx} \mathcal{F}, \quad (4.81)$$

and

$$\text{cov}(\mathbf{x}, \mathbf{u}) = \text{cov}(\mathbf{F} \mathbf{u}, \mathbf{u}) = \mathbf{F} \text{cov}(\mathbf{u}, \mathbf{u}) = \mathbf{F}. \quad (4.82)$$

Eqs (4.81) and (4.82) imply

$$\mathbf{F} = \mathbf{C}_{xx} \mathcal{F}. \quad (4.83)$$

Similarly  $\mathbf{G} = \mathbf{C}_{yy} \mathcal{G}. \quad (4.84)$

Hence the inverse mappings  $\mathbf{F}$  and  $\mathbf{G}$  (from the canonical variates to  $\mathbf{x}$  and  $\mathbf{y}$ ) can be calculated from the forward mapping  $\mathcal{F}^T$  and  $\mathcal{G}^T$ . The matrix  $\mathbf{F}$  is composed of column vectors  $\mathbf{F}_j$ , and  $\mathbf{G}$ , of column vectors  $\mathbf{G}_j$ .  $\mathbf{F}_j$  and  $\mathbf{G}_j$  are the *canonical correlation patterns* associated with  $u_j$  and  $v_j$ , the canonical variates. In general, orthogonality of vectors within a set is not satisfied by any of the four sets  $\{\mathbf{F}_j\}$ ,  $\{\mathbf{G}_j\}$ ,  $\{\mathbf{f}_j\}$  and  $\{\mathbf{g}_j\}$ , and

$$\text{cov}(u_i, u_j) = \text{cov}(v_i, v_j) = \text{cov}(u_i, v_j) = 0 .$$

Fig.4.7 schematically illustrates the canonical correlation patterns.

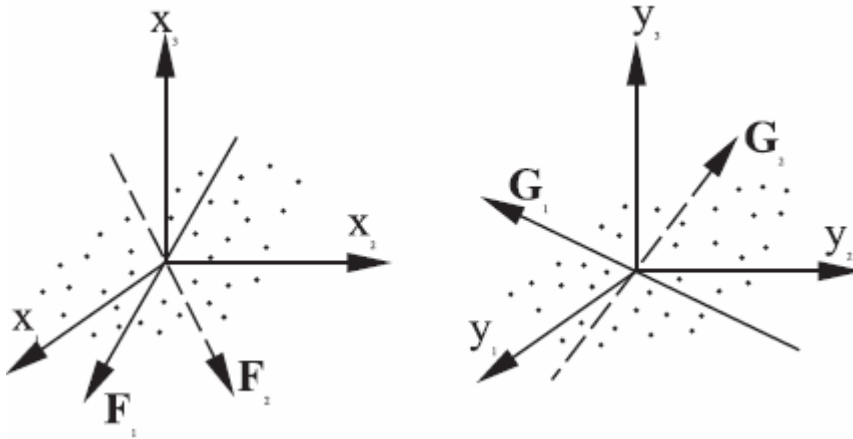


Fig4.7

### 3.2.2 Pre-filter with PCA

When  $x$  and  $y$  contain many variables, it is common to use PCA to pre-filter the data to reduce the dimensions of the datasets, i.e. apply PCA to  $x$  and  $y$  separately, extract the leading PCs, then apply CCA to the leading PCs of  $x$  and  $y$ .

Using Hotelling's choice of scaling for the PCAs (Eq 4.30), we express the PCA expansions as

$$x = \sum_j a'_j e'_j, \quad y = \sum_j a''_j e''_j . \tag{4.85}$$

CCA is then applied to

$$\tilde{x} = [a'_1, \dots, a'_{m_x}]^T, \quad \tilde{y} = [a''_1, \dots, a''_{m_y}]^T , \tag{4.86}$$

where only the first  $m_x$  and  $m_y$  modes are used. Another reason for using the PCA pre-filtering is that when the number of variables is not small relative to the number of samples, the CCA method may become unstable, as the many higher modes may by chance attain high correlation although they account for negligible variance.

With Hotelling's scaling

$$\text{cov}(a'_j, a'_k) = \delta_{jk}, \quad \text{cov}(a''_j, a''_k) = \delta_{jk}, \quad (4.87)$$

leading to

$$\mathbf{C}_{\tilde{x}\tilde{x}} = \mathbf{C}_{\tilde{y}\tilde{y}} = \mathbf{I}. \quad (4.88)$$

Eqs. (4.75) and (4.76) simplify to

$$\mathbf{C}_{\tilde{x}\tilde{y}} \mathbf{C}_{\tilde{x}\tilde{y}}^T \mathbf{f} \equiv \mathbf{M}_f \mathbf{f} = \lambda \mathbf{f}, \quad (4.89)$$

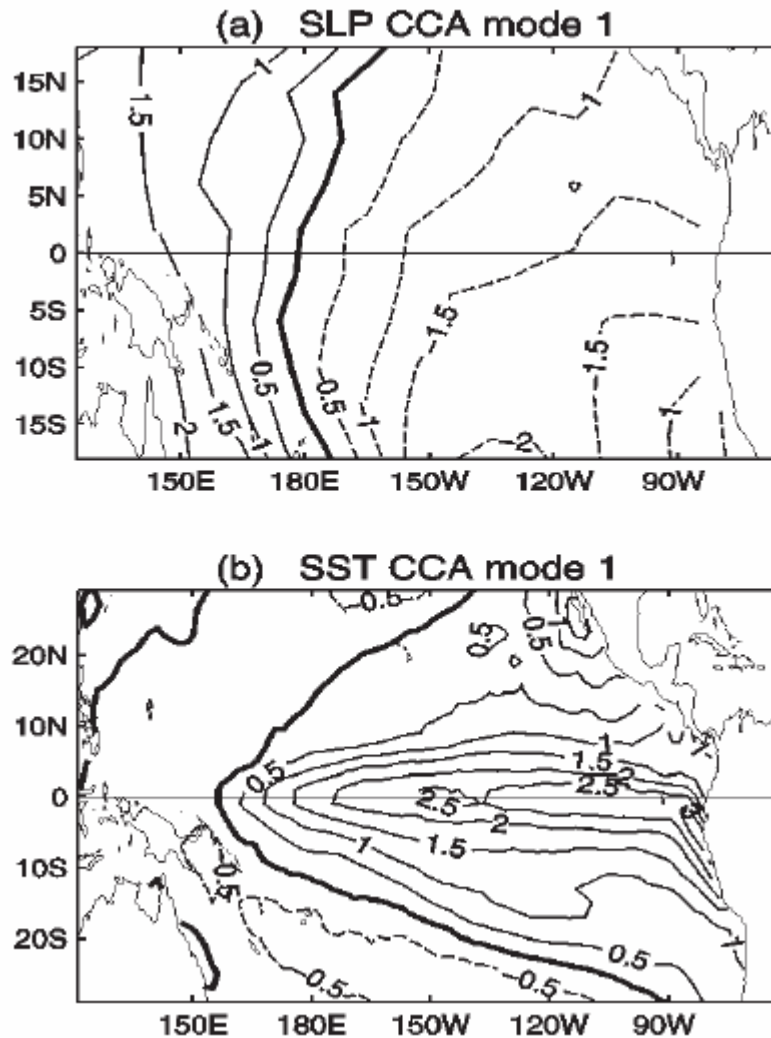
$$\mathbf{C}_{\tilde{x}\tilde{y}}^T \mathbf{C}_{\tilde{x}\tilde{y}} \mathbf{g} \equiv \mathbf{M}_g \mathbf{g} = \lambda \mathbf{g}. \quad (4.90)$$

As  $M_f$  and  $M_g$  are non-negative definite symmetric matrices, the eigenvectors  $\{\mathbf{f}_j\}$   $\{\mathbf{g}_j\}$  are now sets of orthogonal vectors. Eqs (4.83) and (4.84) simplify to

$$\mathbf{F} = \mathcal{F}, \quad \mathbf{G} = \mathcal{G}. \quad (4.91)$$

Hence  $\{\mathbf{F}_j\}$  and  $\{\mathbf{G}_j\}$  are also two sets of orthogonal vectors, and are identical to  $\{\mathbf{f}_j\}$  and  $\{\mathbf{g}_j\}$ , respectively. Because of these nice properties, pre-filtering by PCA (with the Hotelling scaling) is recommended when  $\mathbf{x}$  and  $\mathbf{y}$  have many variables (relative to the number of samples). However, the orthogonality only holds in the reduced dimensional spaces,  $\tilde{\mathbf{x}}$  and  $\tilde{\mathbf{y}}$ . If transformed into the original space  $\mathbf{x}$  and  $\mathbf{y}$ ,  $\{\mathbf{F}_j\}$  and  $\{\mathbf{G}_j\}$  are in general not two sets of orthogonal vectors.

Fig4.8 shows the mode 1 CCA of the tropical Pacific sea level pressure (SLP) field and the SST field, showing clearly the Southern Oscillation pattern in the SLP and the El Niño/La Niña pattern in the SST. The canonical variates  $u$  and  $v$  (not shown) fluctuate with time, both attaining high values during El Niño, low values during La Niña, and neutral values around zero during normal conditions.



### 3.2.3 Singular value decomposition (SVD)

Instead of maximizing the correlation as in CCA, one can maximize the covariance between two datasets. This alternative method is often called the singular value decomposition (SVD).

SVD is identical to CCA except that it maximizes the covariance instead of the correlation. As mentioned before, CCA can be somewhat unstable (when the number of variables is not small relative to the number of samples) in that modes with high correlation may account for negligible variance, hence the recommended pre-filtering of data by PCA before applying CCA. SVD, by using covariance instead of correlation, does not have the unstable nature of the CCA, and does not need the pre-filtering by PCA.



In SVD, one simply performs SVD on the data covariance matrix  $C_{xy}$

$$C_{xy} = USV^T, \quad (4.92)$$

Where the matrix U contains the left singular vectors  $f_i$ , V the right singular vectors  $g_i$ , and S the singular values. Maximum covariance between  $u_i$  and  $v_i$  is attained (Bretherton et al. 1992)

$$u_i = f_i^T \mathbf{x}, \quad v_i = g_i^T \mathbf{y}. \quad (4.93)$$

The inverse transform is given by

$$\mathbf{x} = \sum_i u_i f_i, \quad \mathbf{y} = \sum_i v_i g_i. \quad (4.94)$$

For most application, SVD yields rather similar results to the CCA (with CCA pre-filtering).

The matrix technique SVD can also be used to solve the CCA problem. Similar to (4.63), we seek

$$u = \mathbf{f}^T \mathbf{x}, \quad v = \mathbf{g}^T \mathbf{y}. \quad (4.95)$$

such that

$$\text{cov}(u, v) = \text{cov}(\mathbf{f}^T \mathbf{x}, \mathbf{g}^T \mathbf{y}) = \mathbf{f}^T \text{cov}(\mathbf{x}, \mathbf{y}) \mathbf{g} \rightarrow \text{maximum} \quad (4.96)$$

subject to

$$\mathbf{f}^T \mathbf{f} = 1 \quad \text{and} \quad \mathbf{g}^T \mathbf{g} = 1 \quad (4.97)$$

The solution is obtained by using Lagrange multiplier

$$L = \mathbf{f}^T \text{cov}(\mathbf{x}, \mathbf{y}) \mathbf{g} + \alpha(\mathbf{f}^T \mathbf{f} - 1) + \beta(\mathbf{g}^T \mathbf{g} - 1)$$

Similar to CCA (4.72)-(4.76), we have

$$C_{xy} \mathbf{g} + 2\alpha \mathbf{f} = 0$$

$$\mathbf{C}_{xy}^T \mathbf{f} + 2\beta \mathbf{g} = 0 \quad (4.98)$$

$$\begin{aligned} \mathbf{C}_{xy} \mathbf{g} &= s_x \mathbf{f} \\ \mathbf{C}_{xy}^T \mathbf{f} &= s_y \mathbf{g} \end{aligned} \quad (4.99)$$

(4.99) implies that the solution  $\mathbf{g}$  and can be solved by a singular value decomposition. The same solution is obtained by substituting the two equations into each other to obtain

$$\begin{aligned} \mathbf{C}_{xy}^T \mathbf{C}_{xy} \mathbf{g} &= 4\alpha\beta \mathbf{g} = \lambda \mathbf{g} \\ \mathbf{C}_{xy} \mathbf{C}_{xy}^T \mathbf{f} &= 4\alpha\beta \mathbf{f} = \lambda \mathbf{f} \end{aligned} \quad (4.100)$$

(4.50) is very similar to (4.75) or (4.76). So, the solution can be obtained by eigen equations.

So, there are two approaches to perform SVD: (1) simply perform SVD on  $\mathbf{C}_{xy} = \text{cov}(\mathbf{x}, \mathbf{y})$ ; (2) solving eigen equations (4.100)