

PCR, CCA, and other pattern-based regression techniques

Michael K. Tippett

International Research Institute for Climate and Society
The Earth Institute, Columbia University

Statistical Methods in Seasonal Prediction, ICTP
Aug 2-13, 2010

Principal component analysis and regression

Key idea of PCA: *data compression*

- ▶ Many colinear variables are replaced by a few independent variables (PCs).
- ▶ The new variables are optimally chosen to approximate the original variables.

[Assumption: “important” components have large variance]

How is PCA useful in regression problems?

(Hint: Regression has problems with many colinear predictors, prefers a few independent predictors)

Pop quiz: Statistics

What is the variance of $a_1x_1 + a_2x_1$ when x_1 and x_2 are independent with unit variance? No cross terms.

$$\text{Var}(a_1x_1 + a_2x_1) = a^2 + a_2^2$$

Pop quiz: Statistics

What is the variance of $a_1x_1 + a_2x_1$ when x_1 and x_2 are independent with unit variance? No cross terms.

$$\text{Var}(a_1x_1 + a_2x_1) = a^2 + a_2^2$$

Pop quiz: PCA

- ▶ What is the variance of a PC (time-series)?
- ▶ What is the correlation between PCs?

Pop quiz: PCA

- ▶ What is the variance of a PC (time-series)?
- ▶ What is the correlation between PCs?

Pop quiz: Linear regression

- ▶ (Easy) If $y = ax$ is regression between anomalies x and y , and x and y have unit variance, what does a measure? (E.g., x , y PCs of centered data)

$$\begin{aligned} a &= \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\sigma_{xy}}{\sigma_x^2} \\ &= \frac{\sigma_{xy}}{\sigma_x^2} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \frac{\sigma_y}{\sigma_x} = \text{Corr}(x, y) \frac{\sigma_y}{\sigma_x} \end{aligned}$$

Pop quiz: Linear regression

- ▶ (Easy) If $y = ax$ is regression between anomalies x and y , and x and y have unit variance, what does a measure? (E.g., x , y PCs of centered data)

$$\begin{aligned} a &= \frac{\text{Cov}(x, y)}{\text{Var}(x)} = \frac{\sigma_{xy}}{\sigma_x^2} \\ &= \frac{\sigma_{xy}}{\sigma_x^2} = \frac{\sigma_{xy}}{\sigma_x \sigma_y} \frac{\sigma_y}{\sigma_x} = \text{Corr}(x, y) \frac{\sigma_y}{\sigma_x} \end{aligned}$$

Pop quiz: Linear regression

- ▶ (Hard) Linear (invertible) transformations of the data transform the regression coefficients the same way.

$$y = Ax$$

$$y' = Ly, \quad x' = Mx$$

$$y' = Ly = LAx = LAM^{-1}Mx = (LAM^{-1})x'$$

(LAM^{-1}) = regression coefficient matrix between x' and y'

“Regression is regression is regression”

Pop quiz: Linear regression

- ▶ (Hard) Linear (invertible) transformations of the data transform the regression coefficients the same way.

$$y = Ax$$

$$y' = Ly, \quad x' = Mx$$

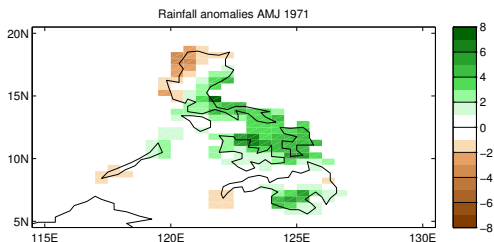
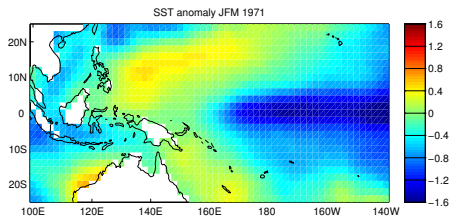
$$y' = Ly = LAx = LAM^{-1}Mx = (LAM^{-1})x'$$

(LAM^{-1}) = regression coefficient matrix between x' and y'

“Regression is regression is regression”

Example: Philippines

Problem: predict gridded April-June precipitation over the Philippines from preceding (January-March) SST.



Example: Philippines

Problem: predict gridded April-June precipitation over the Philippines from preceding (January-March) sea surface temperature.

Details:

- ▶ Data from 1971-2007 (37 years).
- ▶ 194 precipitation gridpoints.
- ▶ 1378 SST gridpoints.

What is the problem?

PCA and regression

For climate forecasts, the length of the historical record severely limits the number of predictors. (Why?)

If the predictors are spatial fields such as SST or the output of a GCM, the number of grid point values (100's, 1000's) is large compared to the number time samples (10's for climate)

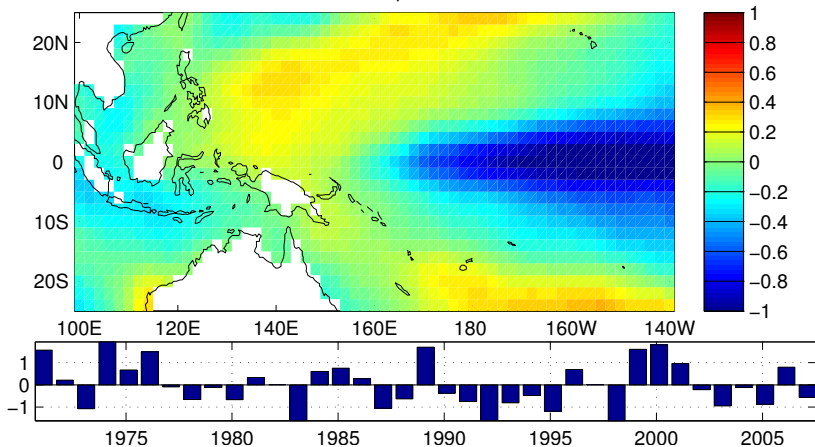
Need to represent the information in the predictor spatial field using fewer numbers.

- ▶ Spatial averages e.g., NINO 3.4.
- ▶ Principal component analysis (PCA).
 - ▶ Weighted spatial average.
 - ▶ Weights are chosen in an optimal manner to maximize explained variance.

Example: PCA of SST

EOF 1 – Correlation with NINO 3.4 = -0.96

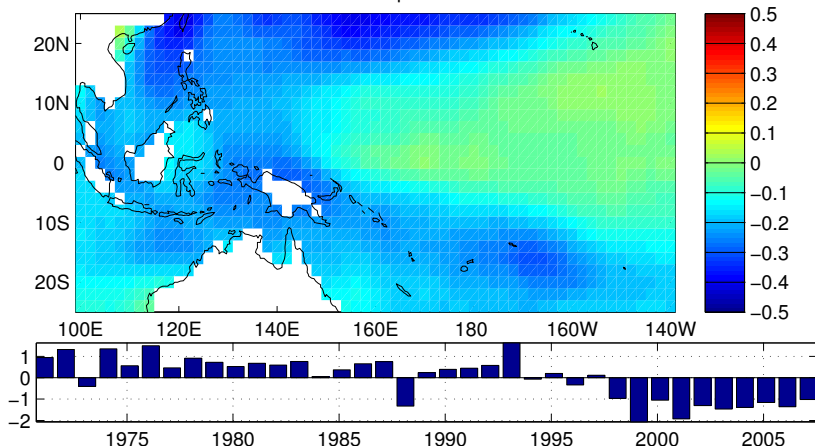
EOF 1 variance explained = 50%



Example: PCA of SST

EOF 2

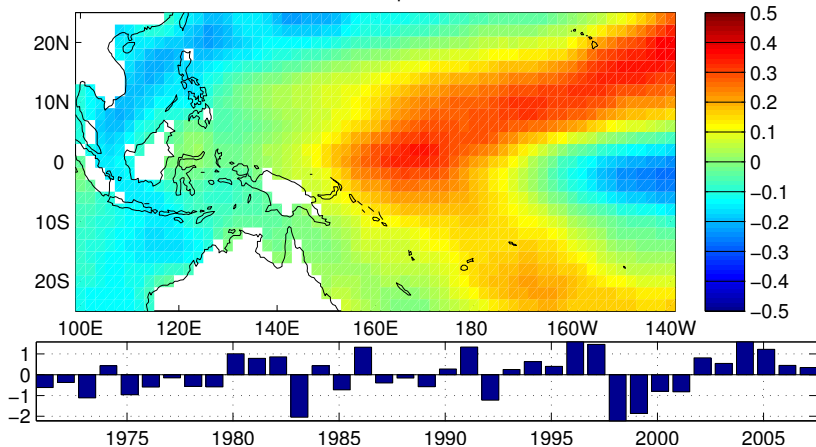
EOF 2 variance explained = 16%



Example: PCA of SST

EOF 3

EOF 3 variance explained = 11%



Principal component regression

PCR

- ▶ $\hat{y} = a_1x_1 + a_2x_2 + \dots a_mx_m + b$
- ▶ Predictors x_i are PCs.

In this example:

- ▶ y = observed precipitation at a gridpoint.
- ▶ PCs of SST anomalies.

How many PCs to use?

Predictor selection problem.

Principal component regression

PCR

- ▶ $\hat{y} = a_1x_1 + a_2x_2 + \dots a_mx_m + b$
- ▶ Predictors x_i are PCs.

In this example:

- ▶ y = observed precipitation at a gridpoint.
- ▶ PCs of SST anomalies.

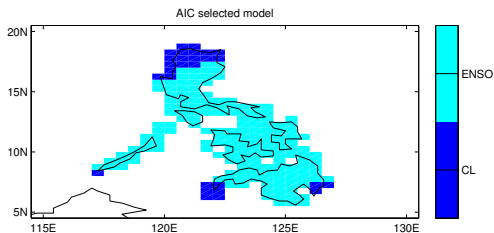
How many PCs to use?

Predictor selection problem.

Example: Philippines

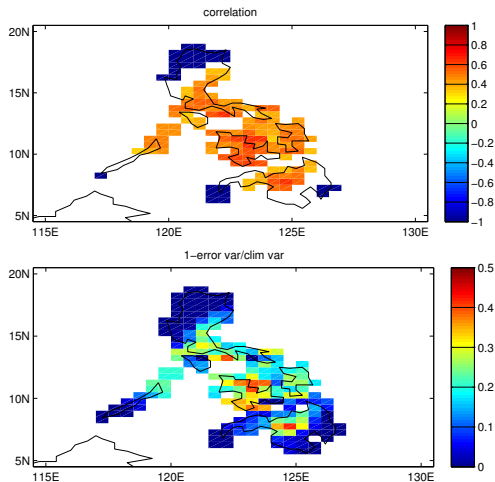
Two models: climatology or ENSO PC as predictor.

For instance, use AIC to select model.



Example: Philippines

- ▶ This model seems to have some skill (cross-validated)
- ▶ Why the negative correlation? [Later]
- ▶ How are the two skill measures related in-sample?



PCA and regression

Is there any benefit to using PCA on the predictand as well as the predictors?

Is there any benefit to predicting the PCs of y rather than y ?

Perhaps. One could imagine a spatial average (like a PC) being more predictable than a value at a gridpoint.

Perhaps not. One could imagine a predictable gridpoint where the PCs explain little variance.

PCA and regression

Is there any benefit to using PCA on the predictand as well as the predictors?

Is there any benefit to predicting the PCs of y rather than y ?

Perhaps. One could imagine a spatial average (like a PC) being more predictable than a value at a gridpoint.

Perhaps not. One could imagine a predictable gridpoint where the PCs explain little variance.

PCA and regression

Is there any benefit to using PCA on the predictand as well as the predictors?

Is there any benefit to predicting the PCs of y rather than y ?

Perhaps. One could imagine a spatial average (like a PC) being more predictable than a value at a gridpoint.

Perhaps not. One could imagine a predictable gridpoint where the PCs explain little variance.

PCA and regression

- ▶ Predicting the PCs of y leads to a different predictor selection problem.
- ▶ Before: select a model for each gridpoint?
- ▶ Now: select a model for each PC?

PCA and regression

- ▶ Predicting the PCs of y leads to a different predictor selection problem.
- ▶ Before: select a model for each gridpoint?
- ▶ Now: select a model for each PC?

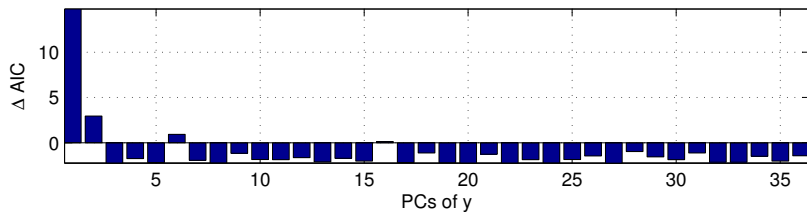
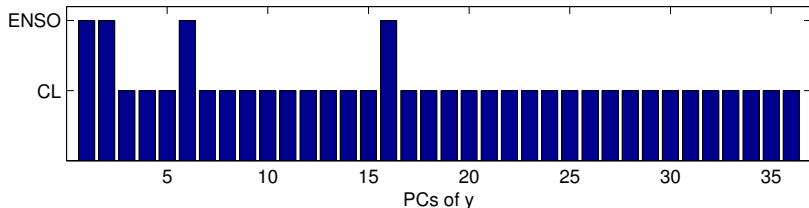
PCA and regression

- ▶ Predicting the PCs of y leads to a different predictor selection problem.
- ▶ Before: select a model for each gridpoint?
- ▶ Now: select a model for each PC?

Example: Philippines

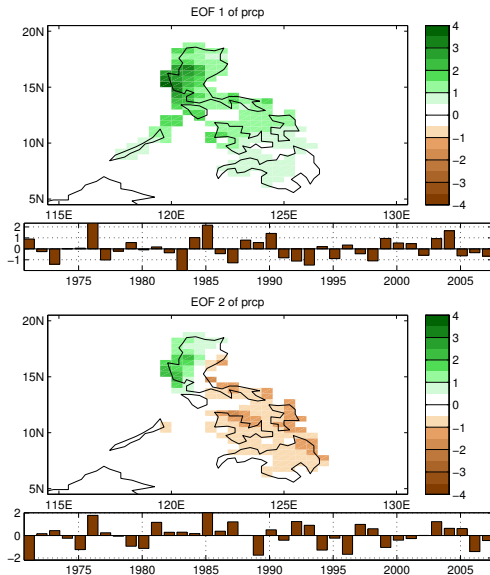
36 PCs of y . (Why?)

Use AIC to select model. ENSO or climatology.



Example: Philippines

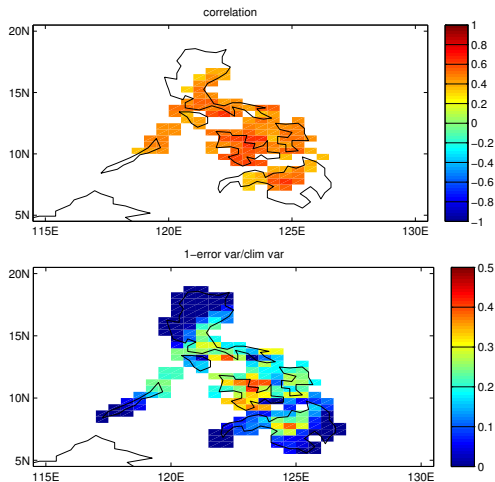
First 2 EOFs of AMJ precipitation:



Example: Philippines

Predicting 2 PCs of rainfall from 1 PC of SST.

- ▶ Correlations of gridpoint and pattern regressions are similar.
- ▶ Normalized error of gridpoint and pattern regressions are similar.



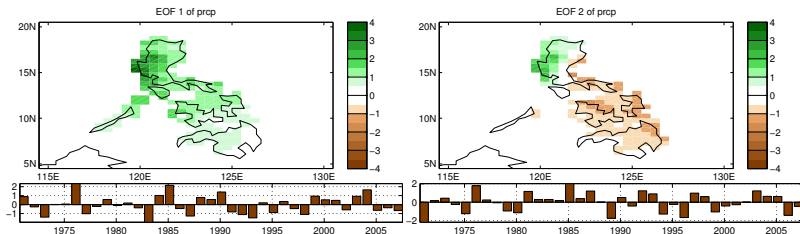
Example: Philippines

Pattern regression is:

$$\hat{PC}_{y1} = 0.61 PC_{x1}$$

$$\hat{PC}_{y2} = -0.36 PC_{x1}$$

What do these numbers mean? (Hint: PCs have unit variance.)



Example: Philippines

Reconstructing the spatial field:

$$\text{Predicted rainfall} = \text{Climatology} + \hat{PC}_{y1} \text{EOF}_{y1} + \hat{PC}_{y2} \text{EOF}_{y2}$$

Difference between prediction and climatology (anomaly) is:

$$\begin{aligned} \text{Predicted anomaly} &= \\ \hat{PC}_{y1} \text{EOF}_{y1} + \hat{PC}_{y2} \text{EOF}_{y2} &= 0.61 \text{PC}_{x1} \text{EOF}_{y1} - 0.36 \text{PC}_{x1} \text{EOF}_{y2} \\ &= \text{PC}_{x1} (0.61 \text{EOF}_{y1} - 0.36 \text{EOF}_{y2}) \end{aligned}$$

Is this simpler? Why?

Predicted anomaly = one time-series \times one pattern.

Example: Philippines

Reconstructing the spatial field:

$$\text{Predicted rainfall} = \text{Climatology} + \hat{P}C_{y1} \text{EOF}_{y1} + \hat{P}C_{y2} \text{EOF}_{y2}$$

Difference between prediction and climatology (anomaly) is:

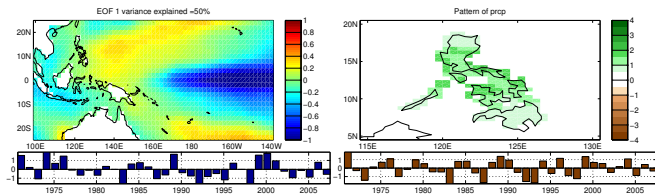
$$\begin{aligned} \text{Predicted anomaly} &= \\ \hat{P}C_{y1} \text{EOF}_{y1} + \hat{P}C_{y2} \text{EOF}_{y2} &= 0.61 \text{PC}_{x1} \text{EOF}_{y1} - 0.36 \text{PC}_{x1} \text{EOF}_{y2} \\ &= \text{PC}_{x1} (0.61 \text{EOF}_{y1} - 0.36 \text{EOF}_{y2}) \end{aligned}$$

Is this simpler? Why?

Predicted anomaly = one time-series \times one pattern.

Example: Philippines

One pattern of rainfall goes with one pattern of SST.



What is the time series of the pattern?

Example: Philippines

Define the pattern to be

$$P \equiv \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ EOF}_{y1} - 0.36 \text{ EOF}_{y2})$$

(Why this scaling???)

The pattern P is a linear combinations of EOFs. What is its time-series TS?

TS is the *same* linear combination of PCs.

The time series TS of the pattern P is:

$$\text{TS} = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ PC}_{y1} - 0.36 \text{ PC}_{y2})$$

What is the variance of TS? (Hint: PCs are independent.) Key

Example: Philippines

Define the pattern to be

$$P \equiv \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ EOF}_{y1} - 0.36 \text{ EOF}_{y2})$$

(Why this scaling???)

The pattern P is a linear combinations of EOFs. What is its time-series TS?

TS is the *same* linear combination of PCs.

The time series TS of the pattern P is:

$$\text{TS} = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ PC}_{y1} - 0.36 \text{ PC}_{y2})$$

What is the variance of TS? (Hint: PCs are independent.) Key

Example: Philippines

Define the pattern to be

$$P \equiv \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ EOF}_{y1} - 0.36 \text{ EOF}_{y2})$$

(Why this scaling???)

The pattern P is a linear combinations of EOFs. What is its time-series TS?

TS is the *same* linear combination of PCs.

The time series TS of the pattern P is:

$$\text{TS} = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ PC}_{y1} - 0.36 \text{ PC}_{y2})$$

What is the variance of TS? (Hint: PCs are independent.) Key

Example: Philippines

Define the pattern to be

$$P \equiv \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ EOF}_{y1} - 0.36 \text{ EOF}_{y2})$$

(Why this scaling???)

The pattern P is a linear combinations of EOFs. What is its time-series TS?

TS is the *same* linear combination of PCs.

The time series TS of the pattern P is:

$$\text{TS} = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ PC}_{y1} - 0.36 \text{ PC}_{y2})$$

What is the variance of TS? (Hint: PCs are independent.) Key

Example: Philippines

Define the pattern to be

$$P \equiv \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ EOF}_{y1} - 0.36 \text{ EOF}_{y2})$$

(Why this scaling???)

The pattern P is a linear combinations of EOFs. What is its time-series TS?

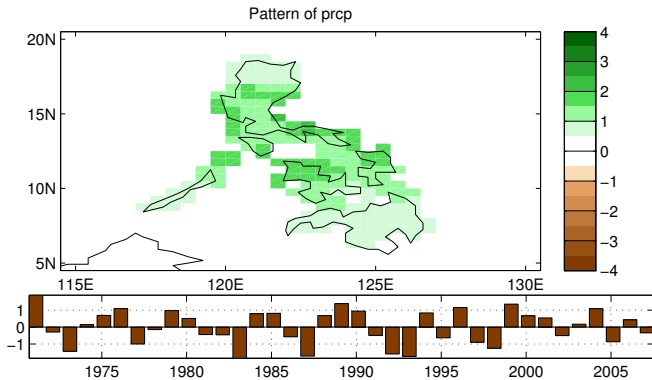
TS is the *same* linear combination of PCs.

The time series TS of the pattern P is:

$$\text{TS} = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 \text{ PC}_{y1} - 0.36 \text{ PC}_{y2})$$

What is the variance of TS? (Hint: PCs are independent.) Key

Pattern P and its time-series TS:



Example: Philippines

The time series TS of the pattern P is:

$$TS = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 PC_{y1} - 0.36 PC_{y2})$$

TS has unit variance.

$$\begin{aligned} \text{Predicted anomaly} &= PC_{x1} (0.61 EOF_{y1} - 0.36 EOF_{y2}) \\ &= PC_{x1} \sqrt{0.61^2 + 0.36^2} P \\ &= 0.71 PC_{x1} P \end{aligned}$$

We are predicting the time-series TS of the pattern P:

$$\hat{TS} = 0.71 PC_{x1}$$

What is 0.71? Hint: TS and PC_{x1} have unit variance.

Example: Philippines

The time series TS of the pattern P is:

$$TS = \frac{1}{\sqrt{0.61^2 + 0.36^2}} (0.61 PC_{y1} - 0.36 PC_{y2})$$

TS has unit variance.

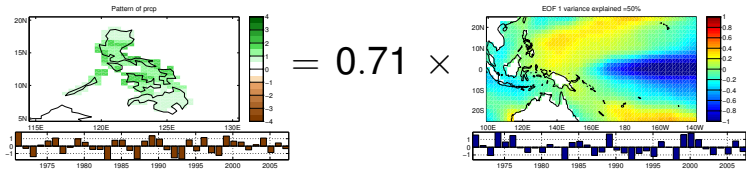
$$\begin{aligned} \text{Predicted anomaly} &= PC_{x1} (0.61 EOF_{y1} - 0.36 EOF_{y2}) \\ &= PC_{x1} \sqrt{0.61^2 + 0.36^2} P \\ &= 0.71 PC_{x1} P \end{aligned}$$

We are predicting the time-series TS of the pattern P:

$$\hat{TS} = 0.71 PC_{x1}$$

What is 0.71? Hint: TS and PC_{x1} have unit variance.

In summary,

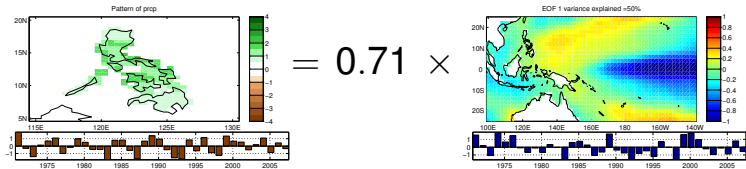


0.71 = correlation between time-series.

In general, *any* pattern regression can be decomposed into *pairs* of patterns related by the *correlation* of their time series.

Canonical correlation analysis (CCA) is an example of such a decomposition.

In summary,



0.71 = correlation between time-series.

In general, *any* pattern regression can be decomposed into *pairs* of patterns related by the *correlation* of their time series.

Canonical correlation analysis (CCA) is an example of such a decomposition.

Pattern regression

$$\begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ \cdot \\ y_l \end{bmatrix} = \begin{bmatrix} a_{11} & a_{12} & \cdot & \cdot & \cdot & a_{1m} \\ a_{21} & a_{22} & \cdot & \cdot & \cdot & a_{2m} \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot & \cdot \\ a_{l1} & a_{l2} & \cdot & \cdot & \cdot & a_{lm} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \cdot \\ \cdot \\ \cdot \\ x_m \end{bmatrix}$$

- ▶ l predictand PCs
- ▶ m predictor PCs
- ▶ $l \times m$ regression coefficients.

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x) \left[\text{Cov}(\text{PC}_x, \text{PC}_x^T) \right]^{-1}$$

Pattern regression

$$y = Ax$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x) \left[\text{Cov}(\text{PC}_x, \text{PC}_x^T) \right]^{-1}$$

What is $\text{Cov}(\text{PC}_x, \text{PC}_x^T)$?

(Hint: PCs are)

$$\text{Cov}(\text{PC}_x, \text{PC}_x^T) = I$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x)$$

What do the elements of A measure? (Hint: PCs are)

$$A = \text{Corr}(\text{PC}_y, \text{PC}_x)$$

Pattern regression

$$y = Ax$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x) \left[\text{Cov}(\text{PC}_x, \text{PC}_x^T) \right]^{-1}$$

What is $\text{Cov}(\text{PC}_x, \text{PC}_x^T)$?

(Hint: PCs are)

$$\text{Cov}(\text{PC}_x, \text{PC}_x^T) = I$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x)$$

What do the elements of A measure? (Hint: PCs are)

$$A = \text{Corr}(\text{PC}_y, \text{PC}_x)$$

Pattern regression

$$y = Ax$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x) \left[\text{Cov}(\text{PC}_x, \text{PC}_x^T) \right]^{-1}$$

What is $\text{Cov}(\text{PC}_x, \text{PC}_x^T)$?

(Hint: PCs are)

$$\text{Cov}(\text{PC}_x, \text{PC}_x^T) = I$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x)$$

What do the elements of A measure? (Hint: PCs are)

$$A = \text{Corr}(\text{PC}_y, \text{PC}_x)$$

Pattern regression

$$y = Ax$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x) \left[\text{Cov}(\text{PC}_x, \text{PC}_x^T) \right]^{-1}$$

What is $\text{Cov}(\text{PC}_x, \text{PC}_x^T)$?

(Hint: PCs are)

$$\text{Cov}(\text{PC}_x, \text{PC}_x^T) = I$$

$$A = \text{Cov}(\text{PC}_y, \text{PC}_x)$$

What do the elements of A measure? (Hint: PCs are)

$$A = \text{Corr}(\text{PC}_y, \text{PC}_x)$$

Pattern regression

$$y = Ax$$

$$A = \begin{bmatrix} \text{Corr}(\text{PC}_{y1}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y1}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y1}, \text{PC}_{xm}) \\ \text{Corr}(\text{PC}_{y2}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y2}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y2}, \text{PC}_{xm}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{Corr}(\text{PC}_{yl}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{yl}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{yl}, \text{PC}_{xm}) \end{bmatrix}$$

In general, each predicted PC of y depends on *all* the PCs of x .

What if A were diagonal? Is it likely that A is diagonal?

correlation = cos angle.

Pattern regression

$$y = Ax$$

$$A = \begin{bmatrix} \text{Corr}(\text{PC}_{y1}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y1}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y1}, \text{PC}_{xm}) \\ \text{Corr}(\text{PC}_{y2}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y2}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y2}, \text{PC}_{xm}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{Corr}(\text{PC}_{yl}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{yl}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{yl}, \text{PC}_{xm}) \end{bmatrix}$$

In general, each predicted PC of y depends on *all* the PCs of x .

What if A were diagonal? Is it likely that A is diagonal?

correlation = cos angle.

Pattern regression

$$y = Ax$$

$$A = \begin{bmatrix} \text{Corr}(\text{PC}_{y1}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y1}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y1}, \text{PC}_{xm}) \\ \text{Corr}(\text{PC}_{y2}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y2}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y2}, \text{PC}_{xm}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{Corr}(\text{PC}_{yl}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{yl}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{yl}, \text{PC}_{xm}) \end{bmatrix}$$

In general, each predicted PC of y depends on *all* the PCs of x .

What if A were diagonal? Is it likely that A is diagonal?

correlation = cos angle.

Pattern regression

$$y = Ax$$

$$A = \begin{bmatrix} \text{Corr}(\text{PC}_{y1}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y1}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y1}, \text{PC}_{xm}) \\ \text{Corr}(\text{PC}_{y2}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{y2}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{y2}, \text{PC}_{xm}) \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots & \vdots & \vdots & \vdots \\ \text{Corr}(\text{PC}_{yl}, \text{PC}_{x1}) & \text{Corr}(\text{PC}_{yl}, \text{PC}_{x2}) & \dots & \dots & \dots & \text{Corr}(\text{PC}_{yl}, \text{PC}_{xm}) \end{bmatrix}$$

In general, each predicted PC of y depends on *all* the PCs of x .

What if A were diagonal? Is it likely that A is diagonal?

correlation = cos angle.

Pattern regression

- ▶ To decompose the regression into pairs of patterns, diagonalize A .
- ▶ Many ways to diagonalize A . The singular value decomposition (SVD) is:

$$A = USV^T$$

where U and S are orthogonal and S is diagonal.
(orthogonal matrix = columns are unit vectors = preserves angles and magnitudes)

Substituting

$$y = Ax = USV^T x$$

or

$$y' = Sx'$$

where

$$y' = U^T y \text{ and } x' = V^T x$$

Pattern regression

- ▶ To decompose the regression into pairs of patterns, diagonalize A .
- ▶ Many ways to diagonalize A . The singular value decomposition (SVD) is:

$$A = USV^T$$

where U and S are orthogonal and S is diagonal.

(orthogonal matrix = columns are unit vectors = preserves angles and magnitudes)

Substituting

$$y = Ax = USV^T x$$

or

$$y' = Sx'$$

where

$$y' = U^T y \text{ and } x' = V^T x$$

Pattern regression

- ▶ To decompose the regression into pairs of patterns, diagonalize A .
- ▶ Many ways to diagonalize A . The singular value decomposition (SVD) is:

$$A = USV^T$$

where U and S are orthogonal and S is diagonal.

(orthogonal matrix = columns are unit vectors = preserves angles and magnitudes)

Substituting

$$y = Ax = USV^T x$$

or

$$y' = Sx'$$

where

$$y' = U^T y \text{ and } x' = V^T x$$

Diagonalized pattern regression

What can we say about the new variables

$$y' = U^T y \text{ and } x' = V^T x?$$

- ▶ y' and x' have unit variance and are uncorrelated (like PCs).
 - ▶ PCs have unit variance and are uncorrelated.
 - ▶ Orthogonal transformation of PCs gives new uncorrelated (angle) variables with unit variance (magnitude).
- ▶ Each new predictand related to just one new predictor.
 - ▶ $y' = Sx'$ where S is diagonal,
- ▶ What do the values of S measure?
(Hint: what is the regression coefficient for two variables with unit variance?)

Diagonalized regression and CCA

This procedure is the same as *canonical correlation analysis*.

Recipe:

- ▶ Regress PCs (uncorrelated unit variance) of y and x .

$$y = Ax$$

- ▶ Use SVD of A to get diagonal relation: $y' = Sx'$.
 - ▶ New variables (canonical variates) are linear (orthogonal) combinations of the PCs.
 - ▶ New variables have unit variance and are uncorrelated.
 - ▶ Associated patterns are linear combinations of EOFs. (Generally not orthogonal).
 - ▶ Elements of S are correlations (canonical correlations).

Diagonalized regression and CCA

This procedure is the same as *canonical correlation analysis*.

Recipe:

- ▶ Regress PCs (uncorrelated unit variance) of y and x .

$$y = Ax$$

- ▶ Use SVD of A to get diagonal relation: $y' = Sx'$.
 - ▶ New variables (canonical variates) are linear (orthogonal) combinations of the PCs.
 - ▶ New variables have unit variance and are uncorrelated.
 - ▶ Associated patterns are linear combinations of EOFs. (Generally not orthogonal).
 - ▶ Elements of S are correlations (canonical correlations).

More CCA

CCA is usually described as finding linear combinations of the x 's and the y 's which have maximum correlation.

Did we do that???

Finding maximum correlation between linear combinations of x and y is the same as finding maximum correlation between linear combinations of x' and y' . Why?

This means we can look at the correlation between linear combinations of x' and y' .

More CCA

CCA is usually described as finding linear combinations of the x 's and the y 's which have maximum correlation.

Did we do that???

Finding maximum correlation between linear combinations of x and y is the same as finding maximum correlation between linear combinations of x' and y' . Why?

This means we can look at the correlation between linear combinations of x' and y' .

More CCA

CCA is usually described as finding linear combinations of the x 's and the y 's which have maximum correlation.

Did we do that???

Finding maximum correlation between linear combinations of x and y is the same as finding maximum correlation between linear combinations of x' and y' . Why?

This means we can look at the correlation between linear combinations of x' and y' .

A calculation

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) = \frac{\text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right)}{\sqrt{\text{Var} \left(\sum a_i x'_i \right) \text{Var} \left(\sum b_j y'_j \right)}}$$

$$\text{Var} \left(\sum a_i x'_i \right) = \sum a_i^2 \text{Var} \left(x'_i \right) = \sum a_i^2 = \|a\|^2$$

$$\text{Var} \left(\sum b_j y'_j \right) = \sum b_j^2 \text{Var} \left(y'_j \right) = \sum b_j^2 = \|b\|^2$$

$$\begin{aligned} \text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right) &= \sum a_i b_i \text{Cov} \left(x'_i, y'_i \right) \\ &= \sum a_i b_i S_i \leq S_1 \sum a_i b_i \leq \|a\| \|b\| \end{aligned}$$

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) \leq S_1$$

A calculation

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) = \frac{\text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right)}{\sqrt{\text{Var} \left(\sum a_i x'_i \right) \text{Var} \left(\sum b_j y'_j \right)}}$$

$$\text{Var} \left(\sum a_i x'_i \right) = \sum a_i^2 \text{Var} \left(x'_i \right) = \sum a_i^2 = \|a\|^2$$

$$\text{Var} \left(\sum b_j y'_j \right) = \sum b_j^2 \text{Var} \left(y'_j \right) = \sum b_j^2 = \|b\|^2$$

$$\begin{aligned} \text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right) &= \sum a_i b_i \text{Cov} \left(x'_i, y'_i \right) \\ &= \sum a_i b_i S_i \leq S_1 \sum a_i b_i \leq \|a\| \|b\| \end{aligned}$$

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) \leq S_1$$

A calculation

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) = \frac{\text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right)}{\sqrt{\text{Var} \left(\sum a_i x'_i \right) \text{Var} \left(\sum b_j y'_j \right)}}$$

$$\text{Var} \left(\sum a_i x'_i \right) = \sum a_i^2 \text{Var} \left(x'_i \right) = \sum a_i^2 = \|a\|^2$$

$$\text{Var} \left(\sum b_j y'_j \right) = \sum b_j^2 \text{Var} \left(y'_j \right) = \sum b_j^2 = \|b\|^2$$

$$\begin{aligned} \text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right) &= \sum a_i b_i \text{Cov} \left(x'_i, y'_i \right) \\ &= \sum a_i b_i S_i \leq S_1 \sum a_i b_i \leq \|a\| \|b\| \end{aligned}$$

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) \leq S_1$$

A calculation

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) = \frac{\text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right)}{\sqrt{\text{Var} \left(\sum a_i x'_i \right) \text{Var} \left(\sum b_j y'_j \right)}}$$

$$\text{Var} \left(\sum a_i x'_i \right) = \sum a_i^2 \text{Var} \left(x'_i \right) = \sum a_i^2 = \|a\|^2$$

$$\text{Var} \left(\sum b_j y'_j \right) = \sum b_j^2 \text{Var} \left(y'_j \right) = \sum b_j^2 = \|b\|^2$$

$$\begin{aligned} \text{Cov} \left(\sum a_i x'_i, \sum b_j y'_j \right) &= \sum a_i b_i \text{Cov} \left(x'_i, y'_i \right) \\ &= \sum a_i b_i S_i \leq S_1 \sum a_i b_i \leq \|a\| \|b\| \end{aligned}$$

$$\text{Corr} \left(\sum a_i x'_i, \sum b_j y'_j \right) \leq S_1$$

More CCA components

Look for the linear combination of x and y that maximizes correlation *but* is uncorrelated with the first component means

- ▶ looking for the linear combination of x'_i and y'_i $i = 2, 3, \dots$ that maximizes correlation. Why?
- ▶ Previous argument give that it is S_2 .

CCA and regression

Knowing CCA is regression is useful ...

- ▶ What happens if many (compared to sample size) PCs are included in a CCA calculation? What happens to the canonical correlations? Overfitting
- ▶ How can the number of PCs included in CCA be decided? Predictor selection, e.g., cross-validation.

CCA and regression

Knowing CCA is regression is useful . . .

- ▶ What happens if many (compared to sample size) PCs are included in a CCA calculation? What happens to the canonical correlations? Overfitting
- ▶ How can the number of PCs included in CCA be decided? Predictor selection, e.g., cross-validation.

CCA and regression

Knowing CCA is regression is useful . . .

- ▶ What happens if many (compared to sample size) PCs are included in a CCA calculation? What happens to the canonical correlations? Overfitting
- ▶ How can the number of PCs included in CCA be decided? Predictor selection, e.g., cross-validation.

CCA and regression

Knowing CCA is regression is useful . . .

- ▶ What happens if many (compared to sample size) PCs are included in a CCA calculation? What happens to the canonical correlations? Overfitting
- ▶ How can the number of PCs included in CCA be decided?
Predictor selection, e.g., cross-validation.

CCA and regression

Knowing CCA is regression is useful . . .

- ▶ What happens if many (compared to sample size) PCs are included in a CCA calculation? What happens to the canonical correlations? Overfitting
- ▶ How can the number of PCs included in CCA be decided? Predictor selection, e.g., cross-validation.

Other pattern regression methods

Other diagonalizations of the regression coefficient matrix A (based on variants of the SVD) give other diagonalized pattern regressions with components that optimize other quantities.

E.g., Redundancy analysis give components that maximize *explained variance*.

Maximum covariance analysis (MCA) finds components with maximum covariance. *However*, the regression between these patterns is generally not diagonal—no simple relations between pairs of patterns.

Summary

- ▶ PCA compresses data and is useful in regressions.
 - ▶ Many \rightarrow few.
 - ▶ Correlated \rightarrow independent
- ▶ In PCR, PCs are the predictors.
- ▶ It can be useful to use PCs as predictands, too.
- ▶ Diagonalizing regressions between PCs decomposes the regression in pairs of patterns.
- ▶ CCA diagonalizes the regression and find the components with maximum correlation.