

Authors

Francesca Elisa Leonelli^{1,3}, Camillo Cammarota³, Salvatore Marullo², Bruno Buongiorno Nardelli¹, Andrea Pisano¹, Chunxue Yang¹

francesca.leonelli@artov.ismar.cnr.it

1 - Consiglio Nazionale delle Ricerche, Istituto di Scienze Marine, CNR-ISMAR, Rome, Italy
2 - Agenzia nazionale per le nuove tecnologie, l'energia e lo sviluppo economico sostenibile (ENEA)
3 - Department of Mathematics, University of Rome "La Sapienza", Rome Italy

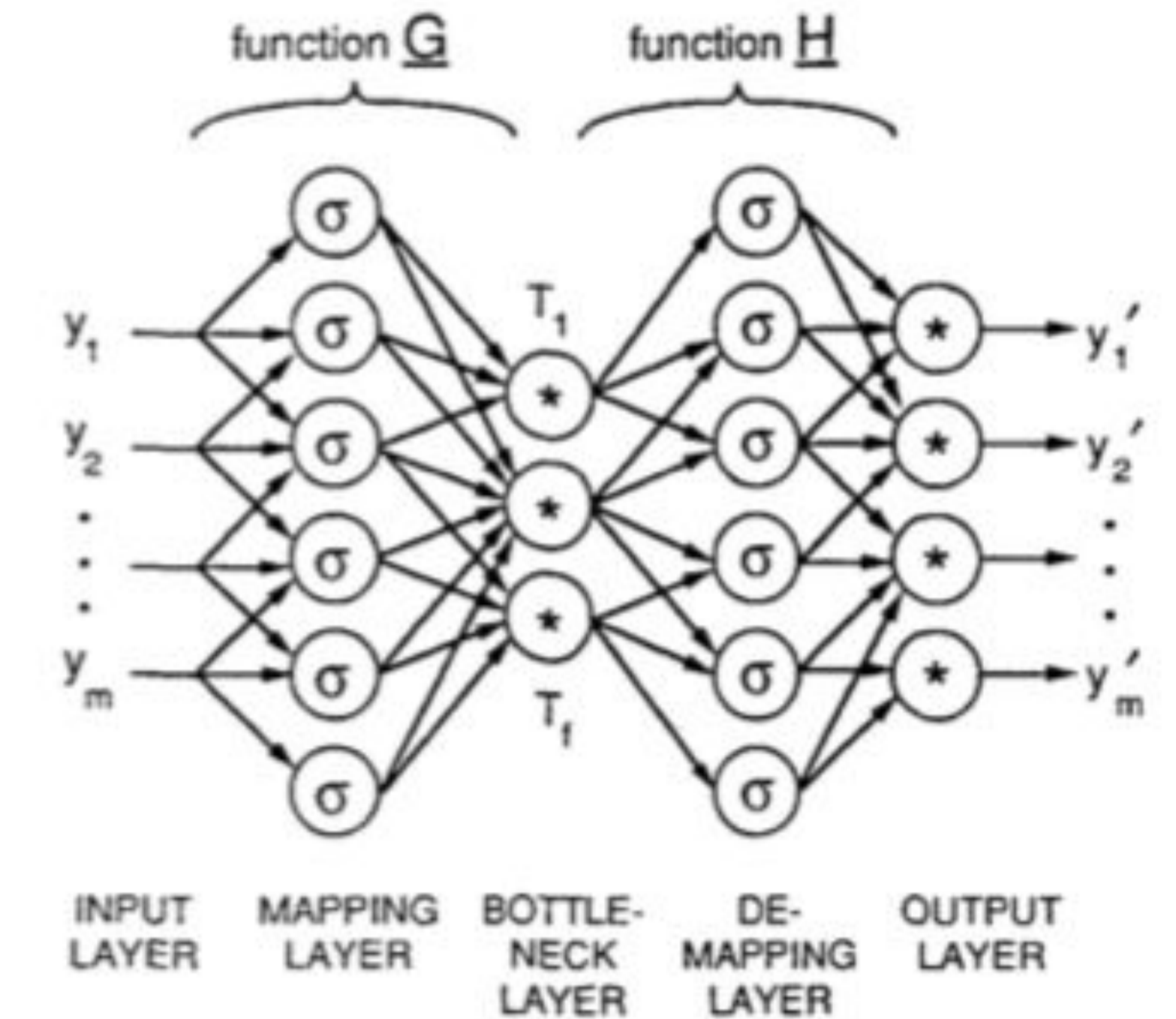


Agenzia nazionale per le nuove tecnologie, l'energia e lo sviluppo economico sostenibile

Context and Motivation

In the environmental sciences, researchers deal with large amounts of data from satellite images of the Earth's surface and from outputs of numerical models. Multivariate techniques for extracting essential information and therefore achieve dimensionality reduction, such as Principal Component Analysis (PCA), become indispensable for treating these datasets. Classical techniques, however, have the strong limitation of finding only linear relations, which does not permit capturing non linear features present in the data, with the risk of misinterpreting these connections and scatter the information into more linear modes. In addition, it is not possible through linear methods to achieve both explaining maximum variance of dataset and capturing local clusters. These limitations can be overcome by introducing neural network models which generalize PCA to a nonlinear multivariate technique of feature extraction (Non Linear PCA). This step is crucial to better describe nonlinear processes typically present in environmental datasets.

Our study is aimed to investigate NLPCA methods on SST products, to evaluate and quantify the improvement in terms of variance explained by the nonlinear modes extracted. In particular we focus on Tropical climate variability, namely El Nino Southern Oscillation (ENSO), Indian Ocean Dipole (IOD), planning to study in detail all Tropical Pacific variability patterns.



Network architecture for simultaneous determination of f nonlinear factors using an autoassociative network [Kramer 1991]

Methods and Results

We consider a shallow feed forward neural network equipped with an *encoder* model (consisting of a layer of input neurons, a hidden encoding layer, and a bottleneck layer) and a *decoder* model (consisting of the encoder's bottleneck layer, a hidden decoding layer and an output layer). The full architecture, that is called an *autoencoder*, can be trained in a supervised scheme, since its objective is to reproduce the input data, by minimizing the loss of information or equivalently maximizing the variance explained.

Given enough hidden neurons, any nonlinear function $f(y)=u$ can be approximated to an arbitrary degree of precision [Cybenko 1989]. After training and validation the neurons of the bottleneck layer will characterize the nonlinear modes retained, accomplishing nonlinear feature extraction, by minimizing the norm of loss of information. (Cost function J minimizes the loss of information, also adopting normalization condition of the bottleneck neuron u , and penalty term on weights is considered to avoid overfitting)

$$u_k = \sum_{j=1}^H w_{jk}^2 \tanh \left(\sum_{i=1}^H w_{ij}^1 y_i + \theta_j^1 \right) + \theta_k^2$$

$$J = (\|y - \tilde{y}\|)^2 + (u)^2 + ((u^2) - 1)^2 + P \sum_{ki} (W_{ki}^{(y)})^2$$

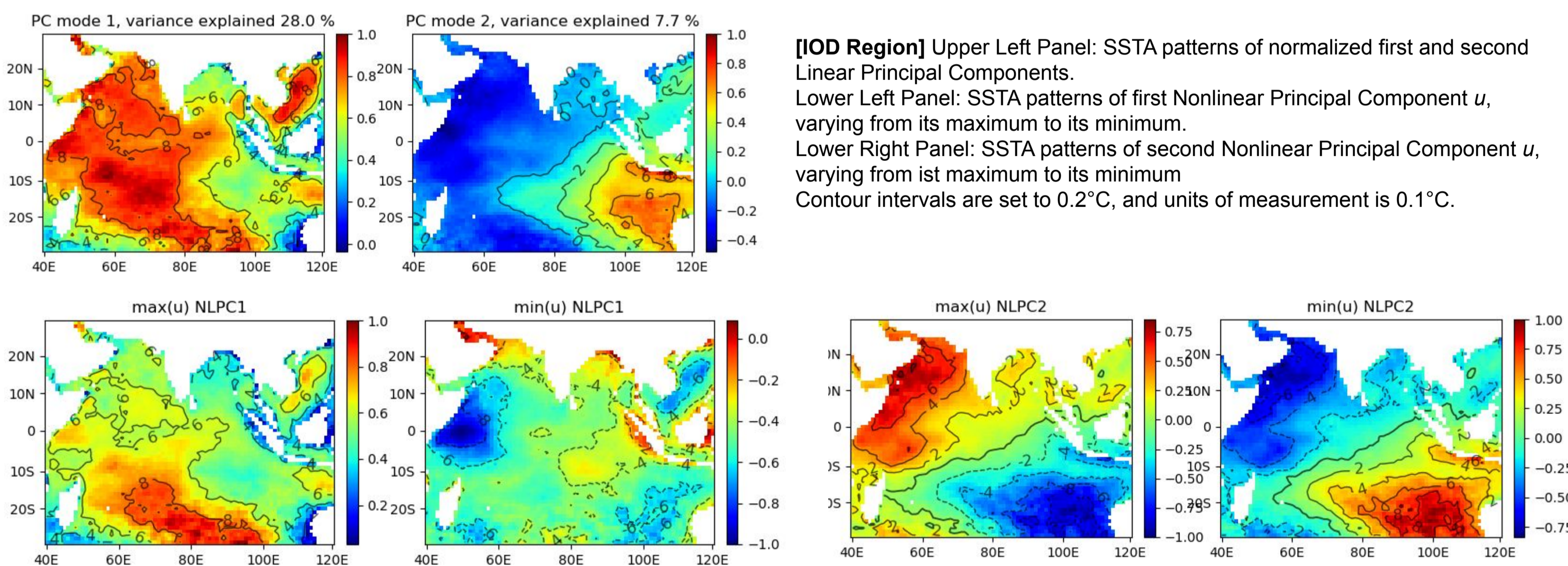
We consider the Sea Surface Temperature ESA CCI SST product with monthly temporal resolution from 1982 to 2018 and $0.25^\circ \times 0.25^\circ$ spatial resolution for the tropical Pacific Ocean, covering the El Nino Southern Oscillation region [30S-30N, 120E-60W], and the Indian Ocean dipole region [30S-30N, 40E-120E]. A climatological annual cycle was calculated in both regions by averaging the data for each calendar month, and monthly SST anomalies (SSTAs) were defined relative to this annual cycle. The SSTA field considered has then been pre-filtered with PCA retaining the first ten components, taking advantage of the data compression aspect of PCA, before using NLPCA for feature extraction.

We therefore perform NLPCA with 10 input-output nodes, 3 mapping-demapping nodes and one node for the bottleneck layer, to retain the first non linear mode. The procedure is repeated on the residual field, retaining a second nonlinear mode.

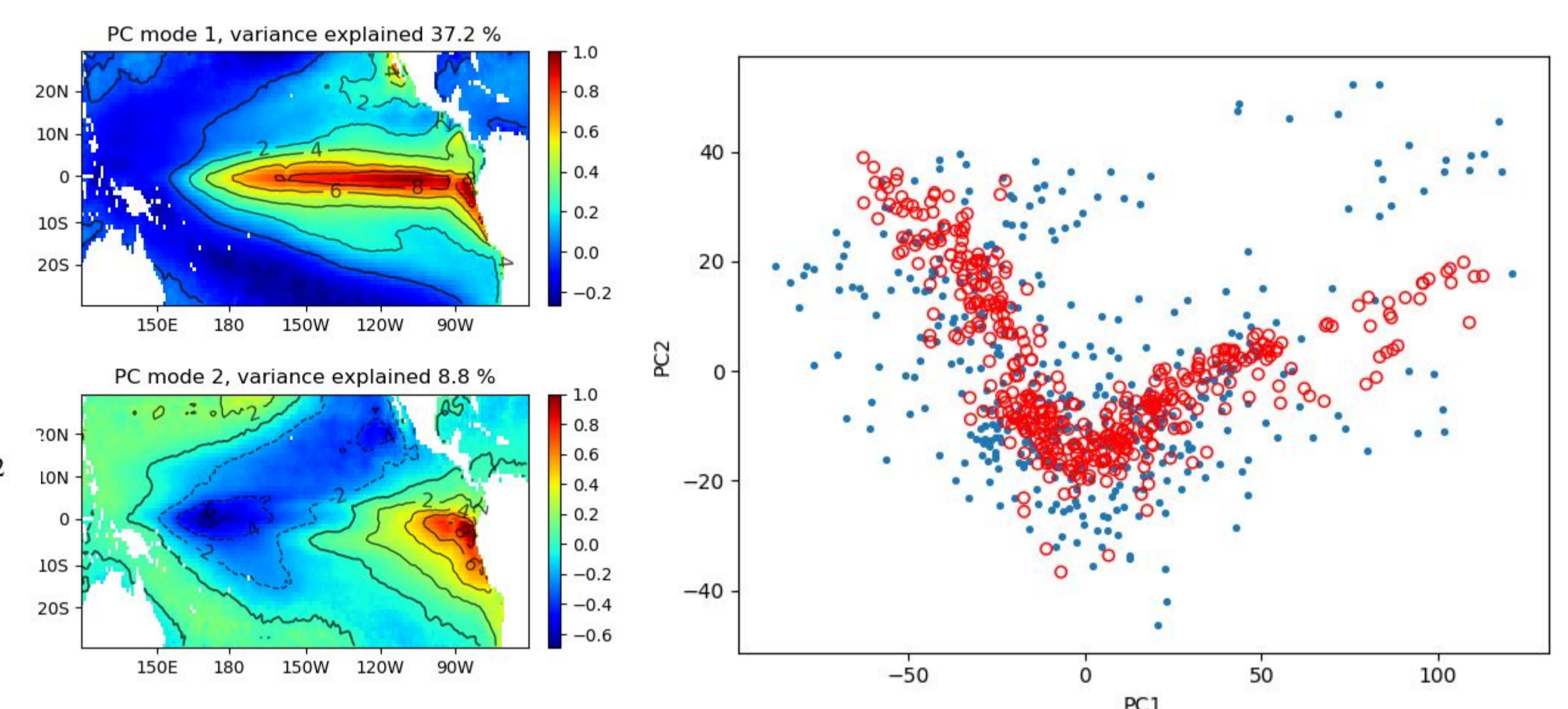
We compare results in terms of variance of the data explained by linear and nonlinear modes:

ENSO	PC mode 1	37.2%	NLPC mode 1	43.1%
	PC mode 2	8.8%	NLPC mode 2	9.4%
IOD	PC mode 1	28.0%	NLPC mode 1	31.6%
	PC mode 2	7.7%	NLPC mode 2	7.3%

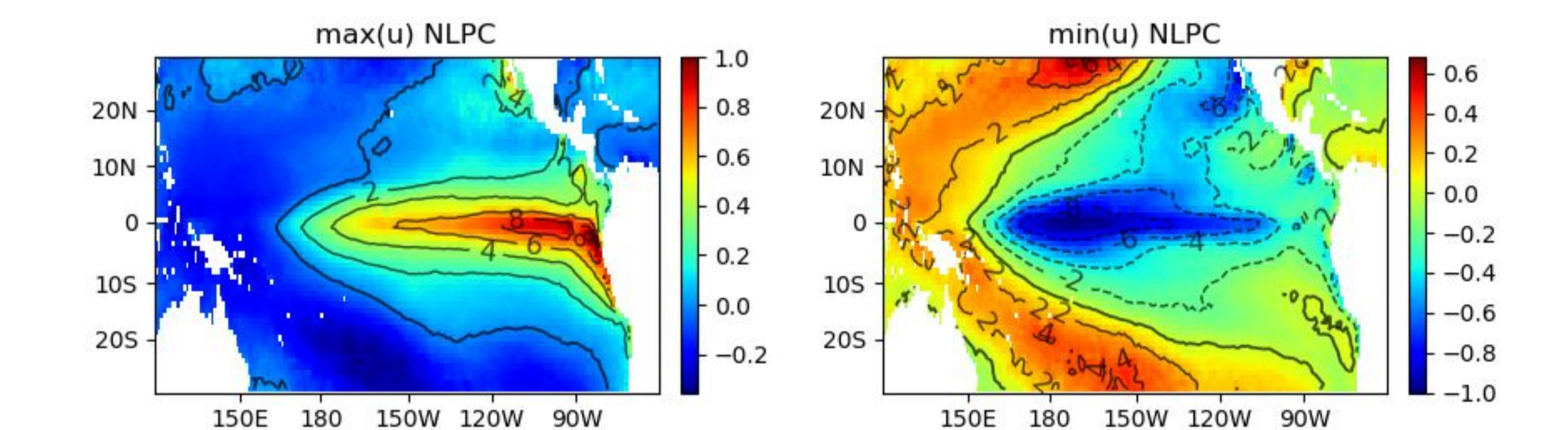
NLPCA gives better results compared to PCA in terms of variance of the data explained by the modes in both regions considered. The lighter discrepancy in IOD region could mean that the oscillatory pattern in this case might be characterized by processes which are more linear compared to the ENSO region.



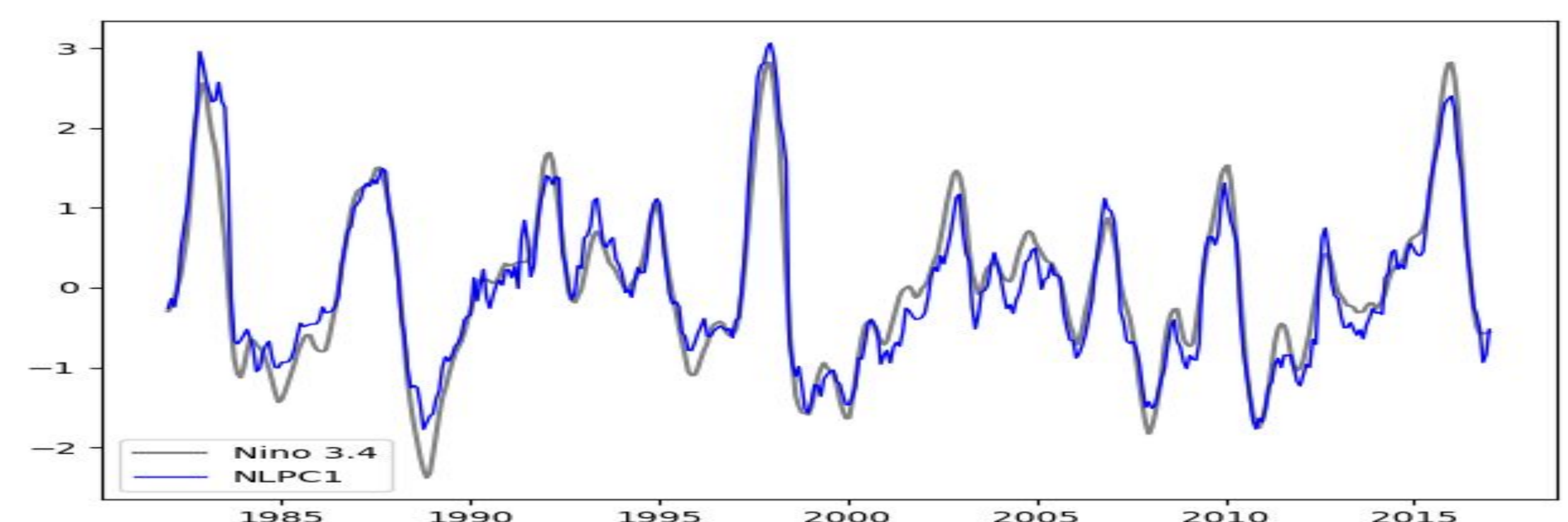
[IOD Region] Upper Left Panel: SSTA patterns of normalized first and second Linear Principal Components. Lower Left Panel: SSTA patterns of first Nonlinear Principal Component u , varying from its maximum to its minimum. Lower Right Panel: SSTA patterns of second Nonlinear Principal Component u , varying from its maximum to its minimum. Contour intervals are set to 0.2°C , and units of measurement is 0.1°C .



[ENSO Region] Left Panel: SSTA patterns of normalized first and second linear Principal Components, contour interval of 0.2°C , units of measurement 0.1°C . Right Panel: Scatterplot of SSTA data in the PC1-PC2 plane, with El Nino states lying in the top right corner and La Nina states lying in the top left corner. First NLPC data is shown in red circles, and captures relation between PC1 and PC2, which clearly are not independent even if uncorrelated.



[ENSO Region] SSTA patterns of first Nonlinear Principal Component u , varying from its maximum (Left panel) to its minimum (Right panel). In contrast to PCA, no single spatial pattern is associated with any given NLPCA mode. The approximation, however, corresponds to a sequence of patterns that can be visualized cinematographically.



[ENSO Region] Plot of normalized time series associated with NLPC mode 1 (blue line) and El Nino 3.4 normalized index (gray line)

Conclusions and Perspectives

- Classical Linear Multivariate techniques such as PCA are powerful tools for the detection of low-dimensional linear structure in multivariate datasets, but can't describe nonlinear processes, typical of environmental data, and tend to scatter a single process into numerous linear modes.
- Neural Network models can accomplish NLPCA, as the natural generalization of the linear feature extraction problem, providing insight of more complex processes with few components.
- A limitation of NLPCA is the non-uniqueness of the solution, which depends on the net architecture choice, on the initialization of the model, and on training/validation parameters selected.
- The ability of the autoencoder can be further investigated by analyzing performance of other types of network layers, typically used for image classification
- Non linear generalization of other classical multivariate techniques, in particular Multichannel Singular Spectrum Analysis, should be further analysed for products' assessment
- Non linear components can be studied to understand if predictability can be better achieved in Tropical Climate Variability

References

- Hsieh, William W. (2004). "Nonlinear multivariate and time series analysis by neural network methods". In: Reviews of Geophysics 42, pp. 10–1029.
- Kramer, Mark A. (1991). "Nonlinear principal component analysis using autoassociative neural networks". In: DOI: 10.1002/aic.690370209.
- Monahan, A. H. (2001). Nonlinear principal component analysis: Tropical Indo-Pacific sea surface temperature and sea level pressure. J. Clim., 14, 219–233.
- Preisendorfer, R. W. (1988). Principal Component Analysis in Meteorology and Oceanography, 425 pp., Elsevier Sci., New York.
- Cybenko, G. 1989. Approximation by superpositions of a sigmoidal function, Mathematics of Control, Signals, and Systems 2, 303–314.