# Spacebased estimation of moisture transport in marine atmosphere using support vector regression

Xiaosu Xie [a,*], W. Timothy Liu [b], Benyang Tang [b]

[a] M/S 300-323, 4800 Oak Grove Dr., Jet Propulsion Laboratory, Pasadena, CA 91109, USA
[b] Jet Propulsion Laboratory, California Institute of Technology, Pasadena, CA, USA

## Abstract

An improved algorithm is developed based on support vector regression (SVR) to estimate horizonal water vapor transport integrated through the depth of the atmosphere ($\Theta$) over the global ocean from observations of surface wind-stress vector by QuikSCAT, cloud drift wind vector derived from the Multi-angle Imaging SpectroRadiometer (MISR) and geostationary satellites, and precipitable water from the Special Sensor Microwave/Imager (SSM/I). The statistical relation is established between the input parameters (the surface wind stress, the 850 mb wind, the precipitable water, time and location) and the target data ($\Theta$ calculated from rawinsondes and reanalysis of numerical weather prediction model). The results are validated with independent daily rawinsonde observations, monthly mean reanalysis data, and through regional water balance. This study clearly demonstrates the improvement of $\Theta$ derived from satellite data using SVR over previous data sets based on linear regression and neural network. The SVR methodology reduces both mean bias and standard deviation compared with rawinsonde observations. It agrees better with observations from synoptic to seasonal time scales, and compare more favorably with the reanalysis data on seasonal variations. Only the SVR result can achieve the water balance over South America. The rationale of the advantage by SVR method and the impact of adding the upper level wind will also be discussed.
© 2007 Elsevier Inc. All rights reserved.

Keywords: Moisture transport; Water cycle; Support vector regression; Remote sensing

## 1. Introduction

Ocean is the main reservoir of heat and water on Earth. The never-ending recycling process in which a small fraction of water is continuously removed from the ocean as excess evaporation over precipitation into the atmosphere, redistributed through atmospheric circulation, deposited as excess precipitation over evaporation on land, and returned to the ocean as river discharge, is critical to the existence of human life and the variability of weather and climate.

The moisture transport integrated over the depth of the atmosphere is

$$\Theta = \frac{1}{g} \int_0^{p_s} q\mathbf{u}\,dp \qquad (1)$$

where $g$ is the acceleration due to gravity, $p$ is the pressure, $p_s$ is the pressure at the surface, $q$ and $\mathbf{u}$ are the specific humidity and wind vector at a certain level. Bold symbols represent vector quantities. $\Theta$ links the water reservoir in the ocean to those over ice and land, and the divergence of which is the difference between evaporation (E) and precipitation (P) at the ocean surface, assuming the change in atmospheric storage is small.

In situ observations for the major components of the hydrological cycle are sparse and intermittent, which severely restrains our understanding of the water cycle. Efforts to retrieve $\Theta$ from rawinsonde observations in various regions and time periods can be traced back several decades ago (e.g., Starr & White, 1955; Rasmusson, 1967; Rosen et al., 1979; Peixoto et al., 1981; Bryan & Oort, 1984). In addition to in situ measurements, numerical weather prediction (NWP) data products have also been used to study atmospheric moisture transport and hydrological budget (e.g., Roads et al., 1992; Trenberth & Guillemot, 1995; Mo & Higgins, 1996; Cohen et al., 2000).

Spacebased observations, with improved coverage and resolution, should improve the estimation of $\Theta$. As suggested

by Liu (1993), $\Theta$ can be written as the product of an equivalent velocity ($\mathbf{u}_e$) and $W$, where

$$W = \frac{1}{g} \int_0^{p_s} q dp \qquad (2)$$

is the precipitable water, and $\mathbf{u}_e = \Theta / W$, by definition, is the depth-average velocity weighted by humidity. Using $W$ derived from the Scanning Multichannel Microwave Radiometer by Liu (1987), two methods to compute monthly mean $\Theta$ over tropical oceans were demonstrated. Heta and Mitsuta (1993) assumed $\mathbf{u}_e$ equals to the 850 mb wind vector produced by the numerical model that assimilated cloud drift wind, while Liu (1993) related $\mathbf{u}_e$ to surface wind vectors derived from satellite data using a polynomial regression. With the availability of surface equivalent neutral wind vector ($\mathbf{u}_s$) from QuikSCAT, Liu and Tang (2005) produced $\Theta$ for tropical and subtropical oceans at twice daily resolution, by relating $\mathbf{u}_e$ to $\mathbf{u}_s$ through neural network (NN). Based on the same physical rationale, Xie and Liu (2005) used nonparametric regression (Hardle, 1990), which is essentially linear regression (LR), to estimate $\Theta$ and showed that the annual mean of $\nabla \cdot \Theta$ bears similar large-scale patterns as that of E-P. Both terms show the major climatological features over the tropical and subtropical oceans.

In this study, an improved method based on SVR (e.g., Vapnik, 1995, 1998; Schölkopf, 1997; Joachims, 1999; among many others), will be presented. The method is described in Section 2 and the training data in Section 3. Validation of the new estimate and the comparison with the linear regression method and the data set produced by Liu and Tang (2005) are shown in Section 4.1 and 4.2. An earlier version of $\Theta$ derived by SVR was used to show continental water balance in South America by Liu et al. (2006); the result is updated in Section 4.3. Discussion is given in Section 5.

## 2. Support vector regression

Support vector machines (SVMs) have several major applications, i.e., classification, regression, and time series prediction. The SVMs for classification have been widely used, such as in handwritten digital recognition, object recognition, speaker identification, face detection in image, and text categorization. The SVMs for regression are referred as support vector regression (SVR) henceforth, which is a statistical tool to derive the relationship between input and output. A comprehensive tutorial of SVMs for pattern recognition can be found in Burges (1998) and a tutorial of SVR can be found in Smola and Schölkopf (2004). In the past decade, SVMs have become increasingly popular due to their good generalization performance. Many studies demonstrated that SVMs match or outperform other machine learning methods (e. g., Blanz et al., 1996; Müller et al., 1997; Drucker et al., 1997; Meyer et al., 2003).

The theoretical basis of SVMs originates from statistical learning theory (Vapnik, 1995, 1998). SVMs essentially translate a nonlinear problem to a linear one by mapping the input into a high-dimensional feature space, and fit a linear model in the feature space (Boser et al., 1992; Vapnik 1995; Cherkassky & Mulier, 1998). The solution of the optimization is unique. Through the introduction of kernel function, the mapping to the

feature space is implicit, and computations are performed directly in the input space rather than in the feature space. The concept of SVR is summarized below through mathematical expressions.

The description here is concise. Details can be referred to Smola and Schölkopf (2004) and Tang and Mazzoni (2006). The training samples are denoted as $(x_i, \Theta_i)\{i = 1, ..., N\}$, where $x_i$ is the input and $\Theta_i$ the output. A kernel mapping function $\Phi(x_i)$ transforms the input into the high-dimensional feature space. Then a linear model is fitted to the data in the feature space by a convex optimization problem:

$$\text{minimize} \frac{1}{2} \| \mathbf{w} \|^2 + C \sum_i (\xi_i + \xi_i^*) \qquad (3)$$

$$\text{subject to} \begin{cases} [\mathbf{w} \cdot \Phi(x_i) + b] - \Theta_i \leq \varepsilon + \xi_i \\ \Theta_i - [\mathbf{w} \cdot \Phi(x_i) + b] \leq \varepsilon + \xi_i^* \\ \xi_i \geq 0, \xi_i^* \geq 0 \end{cases} \qquad (4)$$

where $\mathbf{w}$ is the weight vector in the feature space, $\xi_i$ and $\xi_i^*$ are called slack variables, measuring deviation of the estimated model output from the target data. The output is represented by a linear function in the feature space: $f(x_i) = \mathbf{w} \cdot \Phi(x_i) + b$, where $b$ is a bias term determined during the SVR training. $\varepsilon$ is the precision parameter. The goal of SVR is to find a function $f(x_i)$ with the deviation from the target $\Theta_i$ no larger than $\varepsilon$ for all the training data. In this sense, only those data points with deviations larger than $\varepsilon$ will contribute to the error function (see support vectors defined in the next paragraph). Bigger $\varepsilon$ leads to fewer support vectors. $C$ is the penalty factor, which determines the trade off between the model complexity and the degree to which deviations larger than $\varepsilon$ are tolerated. If $C$ is too large (small), it may cause overfitting (underfitting).

Introducing the Lagrangian multipliers $\alpha_i$ and $\alpha_i^*$, Eqs. (3) and (4) become

$$\text{maximize}: \begin{array}{l} -\frac{1}{2} \sum_{i,j} (\alpha_i^* - \alpha_i)(\alpha_j^* - \alpha_j) \Phi(x_i) \cdot \Phi(x_j) \\ -\varepsilon \sum_i (\alpha_i^* + \alpha_i) + \sum_i \Theta_i(\alpha_i^* - \alpha_i) \end{array} \qquad (5)$$

$$\text{subject to} \sum_{i,j} (\alpha_i^* - \alpha_i) = 0 \text{ and } \alpha_i, \alpha_i^* \in [0, C] \qquad (6)$$

where the weight vector $\mathbf{w}$ is described as a linear combination of the following training pattern

$$\mathbf{w} = \sum_i (\alpha_i^* - \alpha_i) \Phi(x_i) \qquad (7)$$

The inputs with nonzero $\alpha_i^* - \alpha_i$ are called support vectors, which are those inputs whose model outputs $f(x_i)$ differ from the target data $\Theta_i$ by at least $\varepsilon$. The number of support vectors is usually much less than that of the training data samples, which leads to a sparse representation.

By substituting Eq. (7), the output $f(x)$ with an input $x$ can be expressed as

$$f(x) = \mathbf{w} \cdot \Phi(x) + b = \sum_i (\alpha_i^* - \alpha_i) \Phi(x_i) \cdot \Phi(x) + b \qquad (8)$$
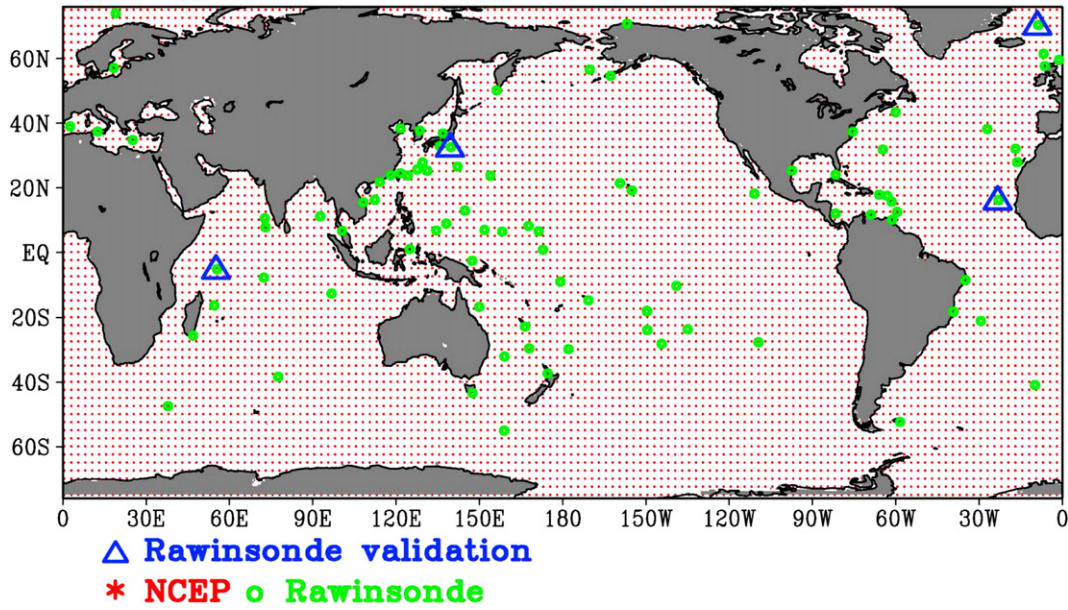
Fig. 1. Distribution of rawinsonde stations, superposed on the NCEP 2.5° grids. The blue triangles mark those rawinsonde stations for validation.

The nonlinear transformation $\Phi(x_i)$ is not required to be explicit. The dot product of the two mapped vectors in Eqs. (5) and (8) is calculated by a kernel function

$$K(x_i, x) = \Phi(x_i) \cdot \Phi(x) \qquad (9)$$

There are many types of kernel functions. In this study, the commonly used radial basis function is adopted

$$K(x_i, x) = \exp\left[-\gamma \, \|x - x_i\|^2\right] \qquad (10)$$

where $\gamma$ is a kernel parameter.

The accuracy of SVR depends on the selection of the hyperparameters (C and $\varepsilon$) and the kernel parameter ($\gamma$). The initial values of the parameters are empirically estimated from the training data based on previous studies. Then only one parameter varies until the optimized correlation between the trained output and the target data is found. In fact, there is a range of the optimized parameters where the correlation coefficient is not sensitive. In this algorithm, $C=10$, $\varepsilon=0.2$, and $\gamma=0.05$.

## 3. Training data

The target data $\Theta$ are daily averages computed from both rawinsonde observations and the National Center for Environmental Prediction (NCEP) reanalysis. The NCEP data are included in the training samples to represent a global coverage, because the rawinsonde measurements are sparse over the ocean
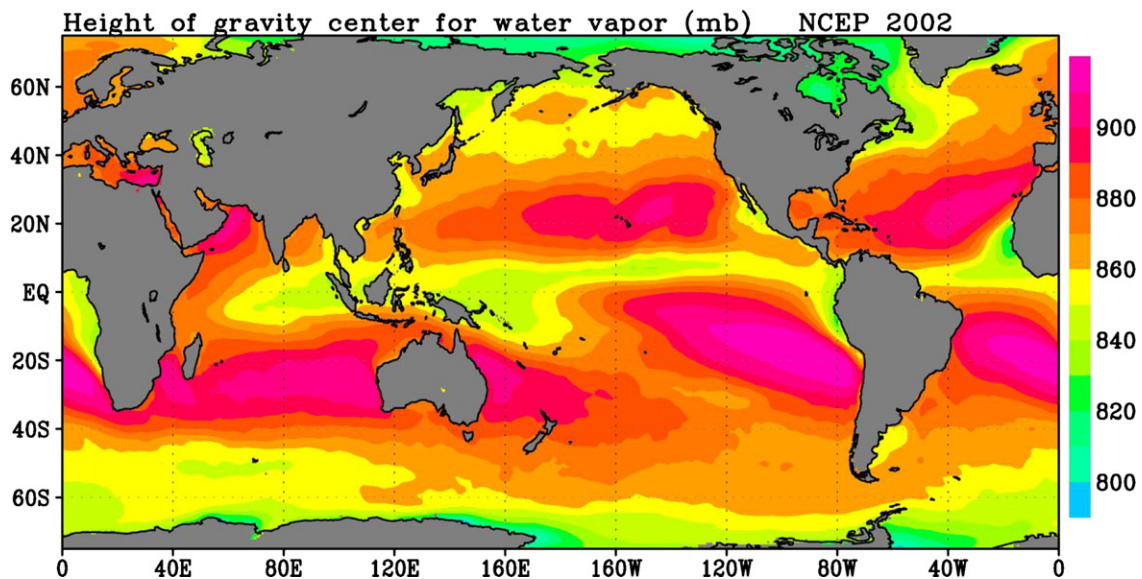


Fig. 2. Pressure of gravity center for humidity estimated from the NCEP reanalysis, averaged for 2002.

Table 1a
Mean and standard deviation of the difference between $\Theta$ from rawinsonde and $\Theta$ from satellite data using the three algorithms

|  | Mean $\Delta\Theta_x$, kg/m/s | Mean $\Delta\Theta_y$, kg/m/s | S.D. $\Delta\Theta_x$, kg/m/s | S.D. $\Delta\Theta_y$, kg/m/s |
|---|---|---|---|---|
| SVM | −2.75 | −8.58 | 69.83 | 60.16 |
| LR | −2.00 | −4.15 | 125.22 | 80.25 |
| NN | −28.56 | −19.00 | 135.19 | 97.54 |

Table 1b
Correlation coefficient between $\Theta$ from rawinsonde and $\Theta$ from satellite data using the three algorithms

|  | $\Delta\Theta_x$, kg/m/s | $\Delta\Theta_y$, kg/m/s |
|---|---|---|
| SVM | 0.948 | 0.867 |
| LR | 0.826 | 0.739 |
| NN | 0.804 | 0.552 |

and with limited geographic coverage (Fig. 1). One year of data in 2001 were used as training samples. 90 rawinsonde stations were selected through data quality control, of which 74 stations are within 40°S–40°N in the three ocean basins and 16 stations are located at higher latitudes. $\Theta$ computed from the NCEP reanalysis is evenly distributed at 2.5° by 2.5° grids between 75°S and 75°N every 10 days. Distribution of the rawinsonde stations is shown in Fig. 1, on top of the NCEP grids. From these data sets, a total of 26,547 data points were selected as the training data set, with 10,000 points randomly selected from the rawinsonde data between 40°S and 40°N, 6547 points poleward of 40° latitude, and 10,000 points randomly selected from the NCEP reanalysis.

The input data include daily averages of $\mathbf{u}_s$ measured by QuikSCAT (Liu, 2002), $\mathbf{u}_{850}$ from the combined cloud drift wind averaged between 800 mb and 900 mb derived from MISR (Horvath & Davies, 2001) and geostationary satellites (Hayden & Pursor, 1995), W from SSM/I (Wentz, 1997), the day of the year, longitude and latitude. Time and longitude take the forms of sine and cosine because of their periodicity. The input parameters are interpolated to the locations of the target data. The input parameters and the target data (x), except for time and longitude, are normalized as: $x' = (x - \bar{x})/\sigma$, where $\bar{x}$ and $\sigma$ are the mean and standard deviation of $x$.
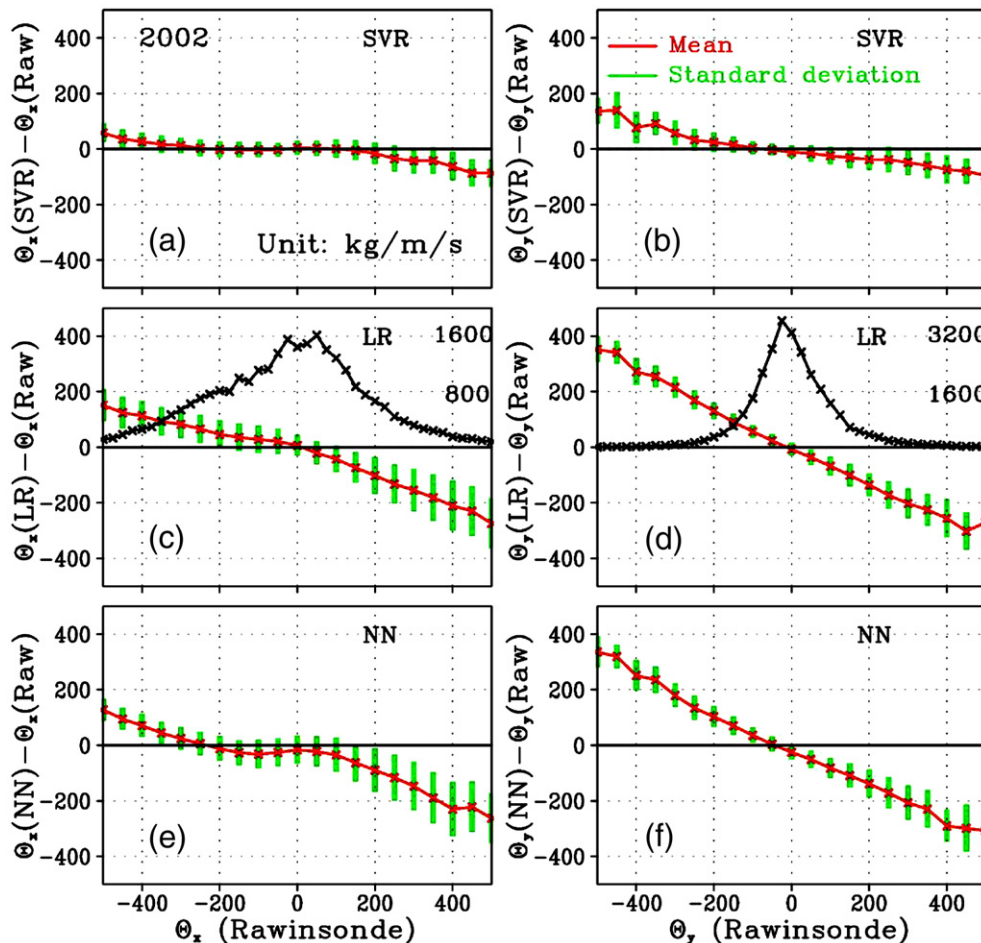


Fig. 3. Bin average of $\Theta_{SVR} - \Theta_R$, as a function of $\Theta_R$, for (a) zonal component, and (b) meridional component. Standard deviation is superimposed on each bin average as error bars. (c) and (d) are the same as (a) and (b), except for $\Theta_{LR} - \Theta_R$. The black curves represent number distributions based on $\Theta_R$. (e) and (f) are the same as (a) and (b), except for $\Theta_{NN} - \Theta_R$.

The methodologies of Liu (1993) and Heta and Mitsuta (1993), as discussed in Section 1, are combined by adding $\mathbf{u}_{850}$ to $\mathbf{u}_s$ as an input parameter. Although humidity peaks near the surface and decreases rapidly with height, the gravity center of the humidity profile is lifted up to the top of the boundary layer over the ascending areas and suppressed to a lower level over the descending regions (Fig. 2). In addition, the dominant mode of vertical humidity variability was found to peak at the top of the boundary layer (Liu et al., 1991), which also suggests the importance of winds at this level in addition to the surface winds.

## 4. Analysis

Four years, 2001–2004, of $\mathbf{\Theta}_{SVR}$ were produced using SVR. The same training data were used to produce $\mathbf{\Theta}_{LR}$, using multivariate linear regression. They are compared with $\mathbf{\Theta}_{NN}$, produced and released to the science community by Liu and Tang (2005), using neural network based on different training data. We did not attempt to optimize the $\mathbf{\Theta}_{NN}$ algorithm in this study. The daily values are first evaluated in Section 4.1 with rawinsonde data in 2002 (independent of the training data), as validation. The distribution of rawinsonde is limited, and monthly mean global distributions were compared with $\mathbf{\Theta}_{NCEP}$ computed from the NCEP reanalysis in Section 4.2.

### 4.1. Comparison with daily rawinsonde data

There are 28,408 rawinsonde observations ($\mathbf{\Theta}_R$) collocated with the satellite-based data. Table 1 shows that the mean difference between $\mathbf{\Theta}_R$ and all the model estimates ($\mathbf{\Theta}_{SVR}$, $\mathbf{\Theta}_{LR}$ and $\mathbf{\Theta}_{NN}$) are small. The standard deviation of the difference is smaller for $\mathbf{\Theta}_{SVR}$ than those for $\mathbf{\Theta}_{LR}$ and $\mathbf{\Theta}_{NN}$. As indicated by the correlation coefficients between model estimates and the rawinsonde observations, $\mathbf{\Theta}_{SVR}$ is better than the other two data sets. The bin-averages in Fig. 3 show that all three satellite-based data sets overestimate compared to observations at low amplitudes and underestimate at high amplitudes, as expected. The numbers of data samples at the two ends are small, as indicated by the $\mathbf{\Theta}_R$ distributions in Fig. 3c and d. The mean difference and standard deviation from the rawinsondes are much lower for $\mathbf{\Theta}_{SVR}$ among the three data sets, particularly for the zonal component.

Performance of each method is also evaluated by time series comparison at individual rawinsonde stations. $\mathbf{\Theta}_{SVR}$ has the best performance across the broad range of time scales. For better figure clarity, $\mathbf{\Theta}_{LR}$ is not shown. Comparisons at two stations, representing different wind regimes, are shown as examples (Figs. 4 and 5). $\mathbf{\Theta}_{SVR}$ captures not only the seasonal variation but also the high-frequency variability of $\mathbf{\Theta}_R$, with amplitude in good agreement with observations. For instance,
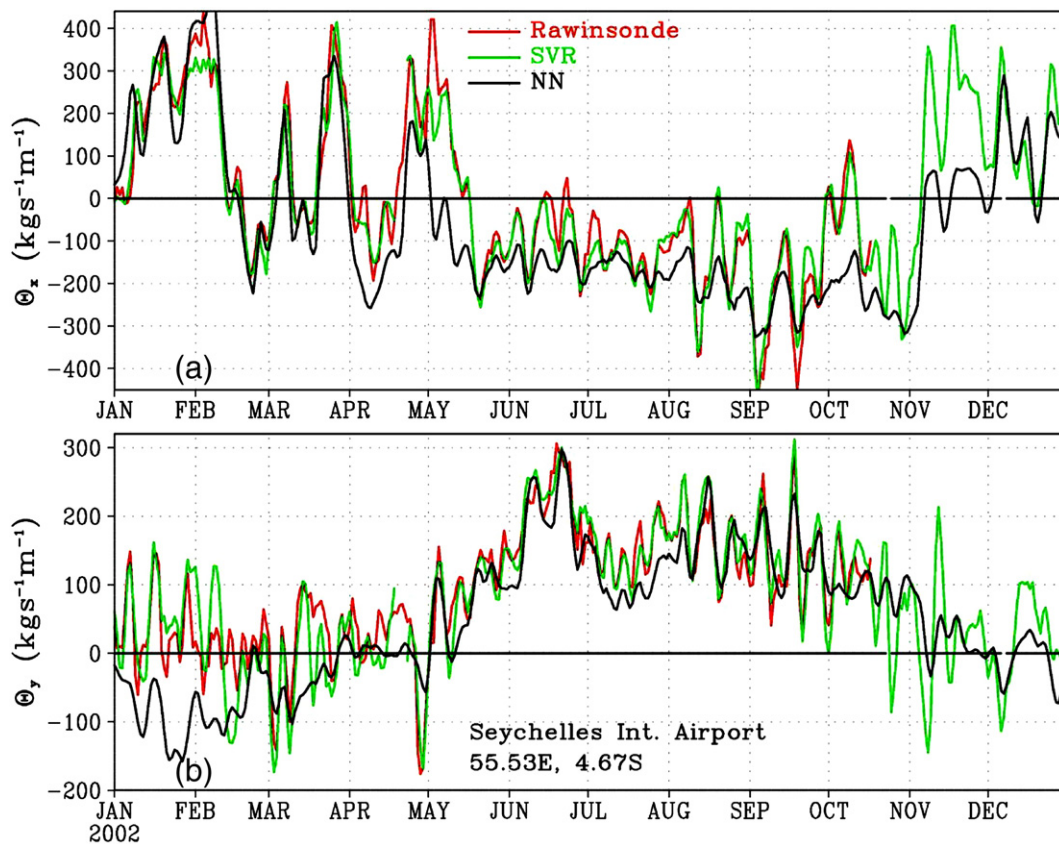


Fig. 4. Comparison of daily $\Theta$ derived from satellite data using SVR (green curve) and NN (black curve) with $\Theta$ derived from rawinsonde (red curve) for (a) zonal and (b) meridional components at Seychelles International Airport (55.53°E, 4.67°S). A 3-day running mean is applied.
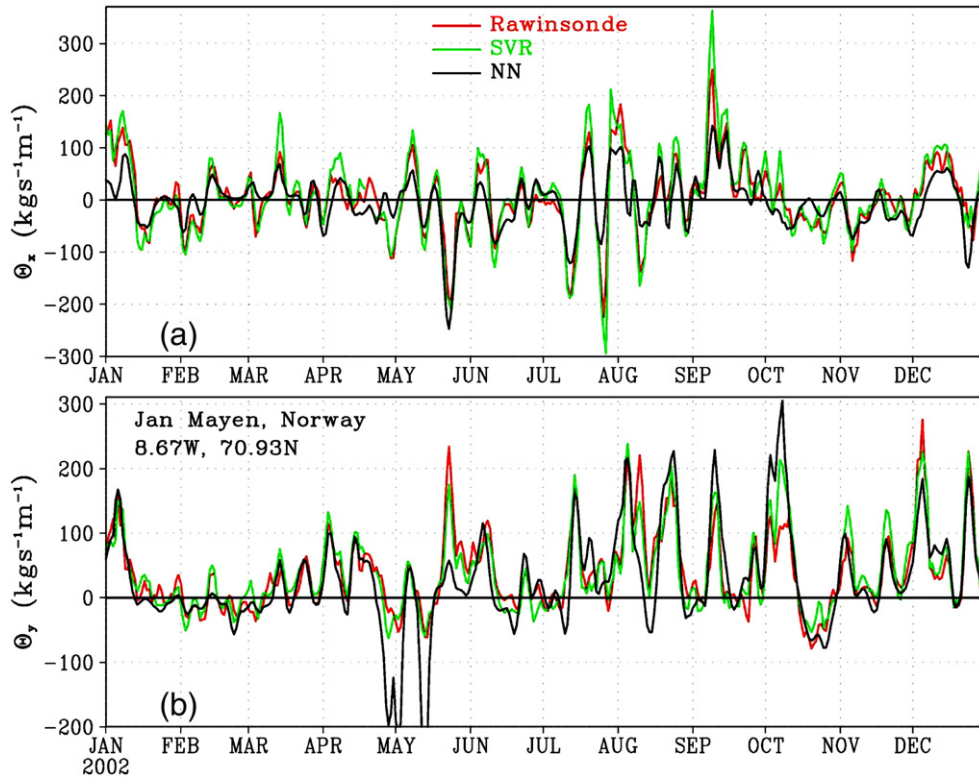
Fig. 5. Same as Fig. 4, except for the rawinsonde station at Jan Mayen, Norway (8.67°W, 70.93°N).
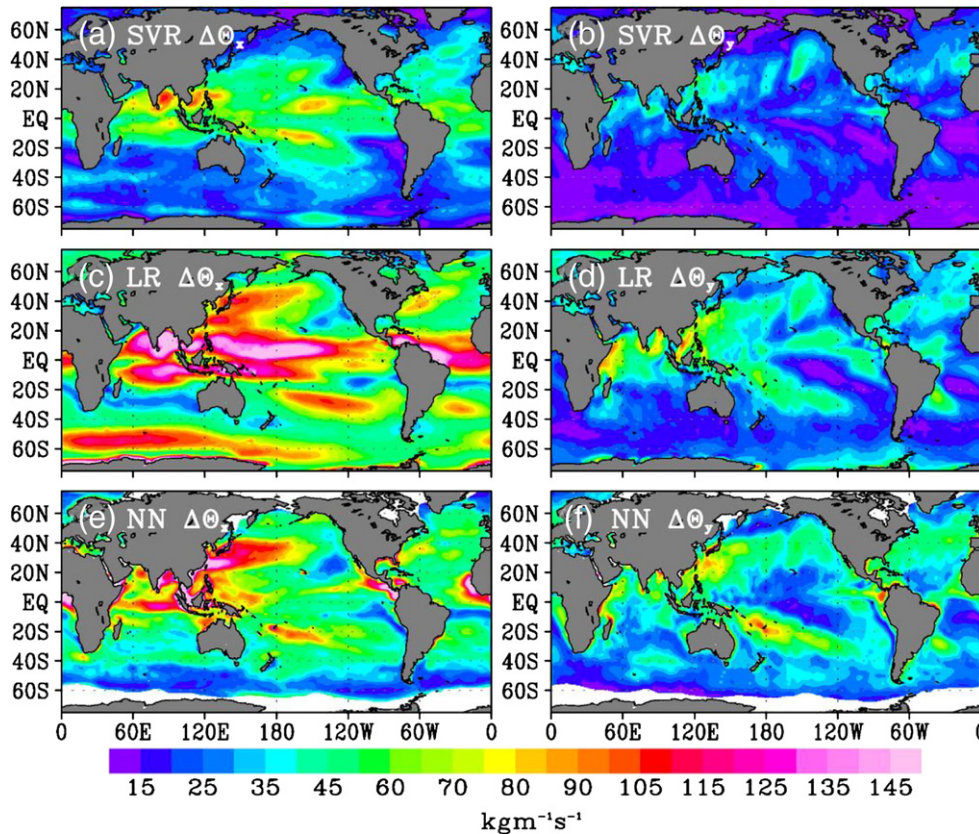


Fig. 6. RMS differences of monthly $\Theta$ derived from satellite data using SVR and from NCEP reanalysis for (a) zonal and (b) meridional components. (c) and (d) are the same as (a) and (b), except for LR. (e) and (f) are the same as (a) and (b), except for NN.

the time series at Seychelles International Airport (Fig. 4) demonstrates the reversal of the Indian summer and winter monsoon. The station at high latitude (Fig. 5 for Jan Mayen, Norway) does not show any clear seasonal cycle but has higher frequency variations. $\Theta_{LR}$ (not shown in the figures) and $\Theta_{NN}$ often lack sufficient skill in capturing the high and low extreme values, which generally produces a narrower range. $\Theta_{LR}$ is particularly poor in high latitude compared with the rawinsonde observations. $\Theta_{NN}$ is sometimes deficient in high-frequency variability, such as the meridional component during January–April 2002 at Seychelles International Airport (Fig. 4). The meridional component of $\Theta_{NN}$ shows unrealistic amplitudes in April and May in 2002 at Jan Mayen, Norway (Fig. 5).

## 4.2. Comparison with monthly NWP estimates

Among the root-mean-square (RMS) differences between the three data sets derived from satellite data and from NCEP reanalysis, the SVR method gives the lowest values in general (Fig. 6). Relatively large RMS in the zonal component of $\Theta_{SVR}$ is found in the Bay of Bengal and South China Sea, where strong seasonal variations of the monsoon circulation dominate. $\Theta_{LR}$ has much larger RMS differences than the other two data sets, especially for the zonal component in the tropics and the southern ocean. Fig. 7 shows that the temporal variation of $\Theta_{NCEP}$ has significant correlation with $\Theta_{SVR}$ over global oceans (with correlation coefficient of over 0.9 in most areas) for the 4-year monthly mean time series. The zonal components of $\Theta_{LR}$ and $\Theta_{NN}$ show low correlation with $\Theta_{NCEP}$ to the west of Japan and to the east along the Kuroshio extension. In this region, the time series in Fig. 8 shows that the low frequency variation of $\Theta_{SVR}$ agrees with $\Theta_R$ marginally better than $\Theta_{NCEP}$, while $\Theta_{LR}$ and $\Theta_{NN}$ fail to capture the peaks of westerly in the summers of 2003 and 2004. The meridional component of $\Theta_{LR}$ and $\Theta_{NN}$ does not agree with $\Theta_{NCEP}$ very well in the temporal variability in the eastern tropical Atlantic off the African coast and the eastern Pacific off Baja California and off Peru (Fig. 7d, f). The zonal component of $\Theta_{LR}$ also shows low correlation with $\Theta_{NCEP}$ in the eastern Pacific and Atlantic, south of the equator, and Gulf Stream. The reason is that $\Theta_{LR}$ poorly predicts the seasonal cycle in these regions; it fails to capture the strong easterlies during the boreal spring in the eastern Pacific and Atlantic, and the strong westerlies during the boreal summer in the Gulf Stream. There is no rawinsonde station in these regions, except in the Gulf Stream. The zonal component of $\Theta$ for the three satellite-based data sets and NCEP reanalysis all capture the seasonal cycle of the rawinsonde observations, but those derived from LR and NN have lower amplitudes. For the meridional components south of Baja California and west of Cape Verde of Africa, all data sets are poor compared to observations, because the amplitudes and the annual cycles are weak, as demonstrated in Fig. 9.
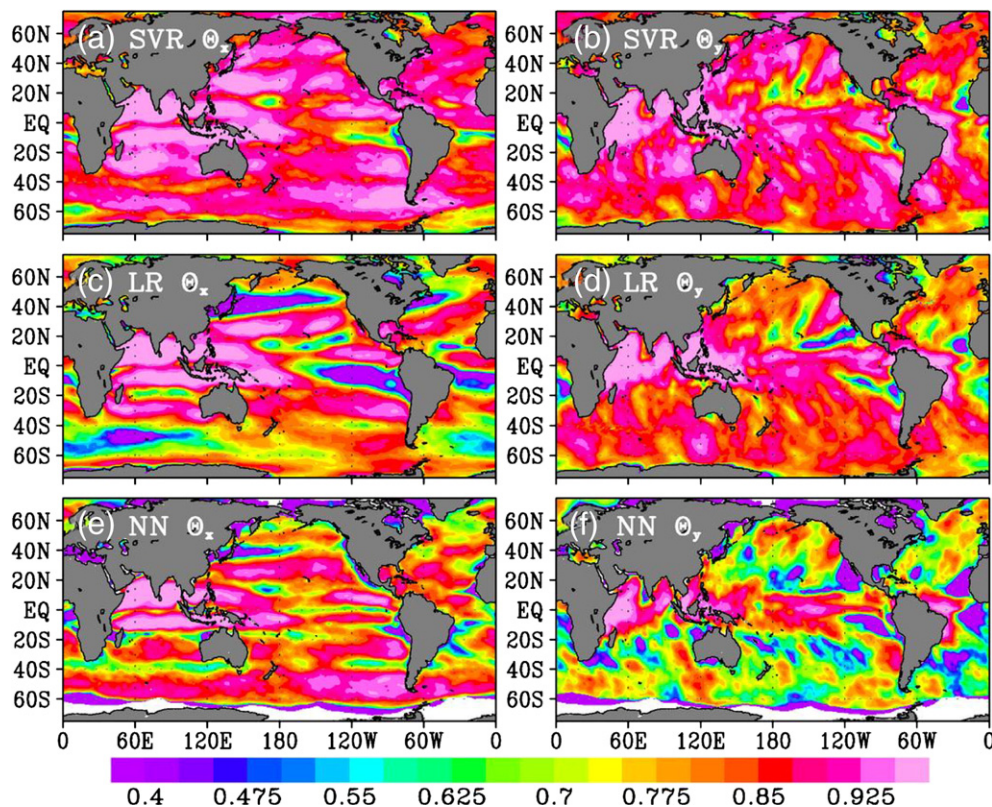


Fig. 7. Correlation coefficients between monthly $\Theta$ derived from satellite data using SVR and from NCEP reanalysis for (a) zonal and (b) meridional components. (c) and (d) are the same as (a) and (b), except for LR. (e) and (f) are the same as (a) and (b), except for NN.
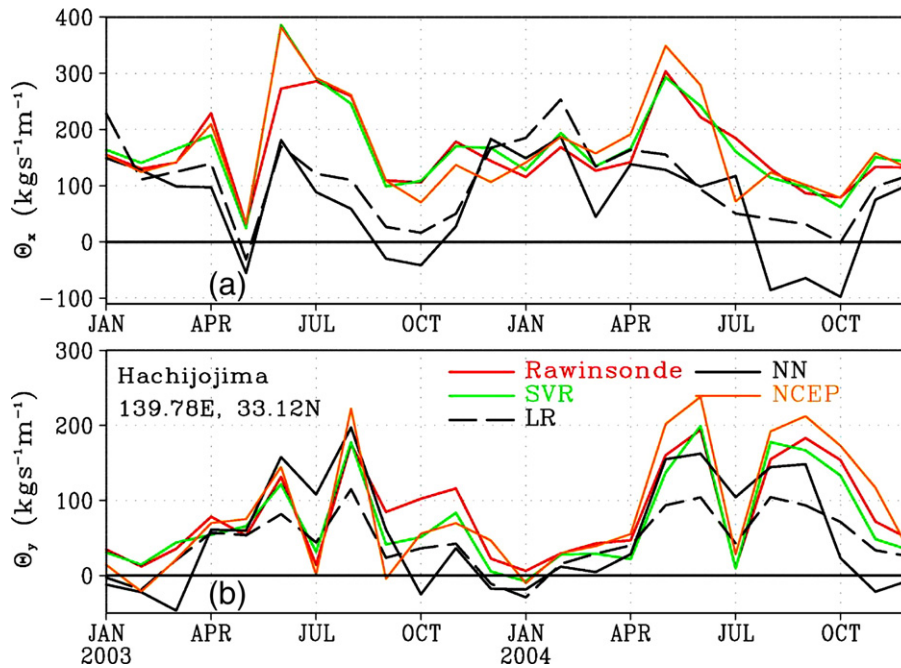
Fig. 8. Comparison of monthly Θ derived from satellite data using SVR (green curve), LR (dashed black curve) and NN (solid black curve) with Θ derived from rawinsonde (red curve) and the NCEP reanalysis (orange curve) for (a) zonal and (b) meridional components at Hachijojima (139.78°E, 33.12°N).

## 4.3. Water balance over South America

The best validation of the remote sensing technique is, perhaps, through application. Liu et al. (2006) showed the approximate balance of the mass change rate ($\partial M/\partial t$) measured by Gravity Recovery and Climate Experiment (GRACE), for the continent of South America, with the normal component of Θ integrated across the entire coast lines of South America ($\int \Theta$)

minus the climatologic river runoff from the continent (R), in agreement with conservation principle.

$$\frac{\partial M}{\partial t} = \int \mathbf{\Theta} - R \tag{11}$$

The Θ used by Liu et al. (2006) is derived by the same SVR technique, except without the use of $\mathbf{u}_{850}$ as input parameter.
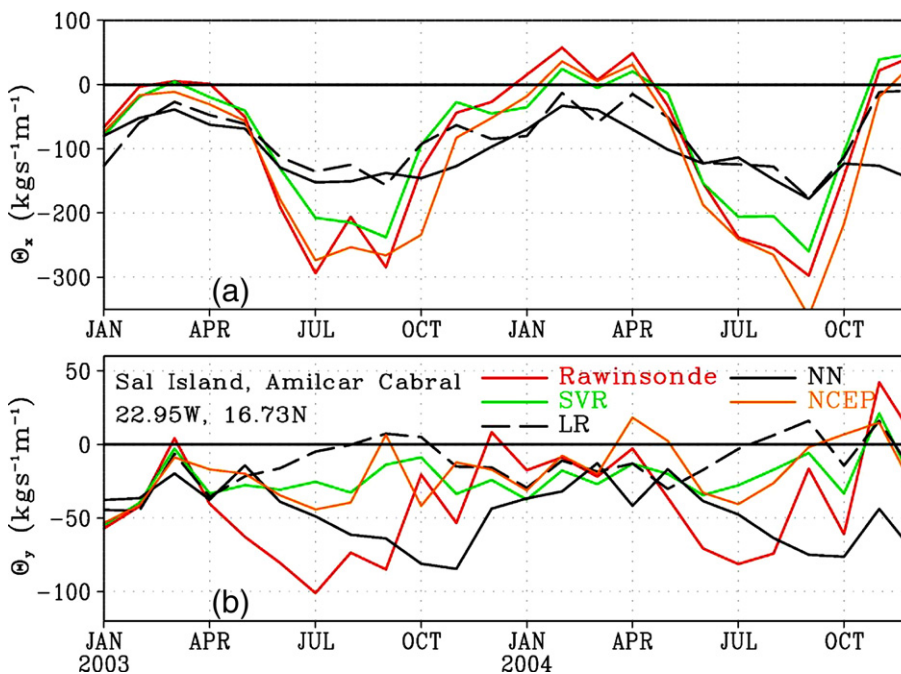


Fig. 9. Same as Fig. 8, except for rawinsonde station at Sal Island, Amilcar Cabral (22.95°W, 16.73°N).
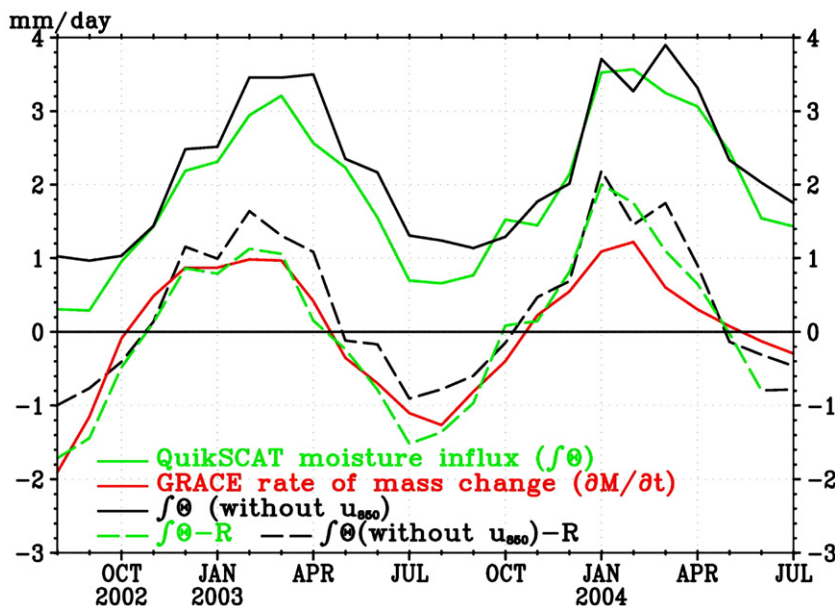
Fig. 10. Annual variation of hydrologic parameters over South America: mass change rate $\partial M/\partial t$ (red line), total moisture transport across coastlines into the continent $\int\Theta$ estimated from satellite data using SVR (solid green line), $\int\Theta - R$ (dashed green line, where $R$ is climatological river discharge), $\int\Theta$ estimated from the earlier version of SVR without $\mathbf{u}_{850}$ (solid black line), and the respective $\int\Theta - R$ (dashed black line).

Using both $\Theta_{LR}$ and $\Theta_{NN}$ result in too much moisture inflow into the continent to balance $\partial M/\partial t$ and $R$. The addition of $\mathbf{u}_{850}$ as input parameter slightly improve the balance as shown in Fig. 10.

## 5. Discussion

In this study, an improved approach using SVR to estimate $\Theta$ over the global ocean from combined satellite observations, including $\mathbf{u}_s$ measured by QuikSCAT, $\mathbf{u}_{850}$ derived from MISR and geostationary satellites, and $W$ from SSM/I, is presented. The statistical model is constructed through SVR by mapping the input parameters, which are $\mathbf{u}_s$, $\mathbf{u}_{850}$, $W$, time and location, to the target data calculated from rawinsondes and the NCEP reanalysis. The results are validated against independent rawinsonde observations and compared with data sets produced based on LR and NN. The SVR algorithm reduces both the mean bias and the standard deviation from the observations, outperforming the other two data sets. It captures not only the seasonal changes but also the synoptic and intraseasonal variations of the observations; the agreement is in both phase and amplitude, and in both low and high latitude oceans. A common deficiency in the previous data sets is their inability to produce the extreme values, and therefore they do not have the full spectrum of variability in the observations.

The monthly mean $\Theta$ derived from the three methods is also compared with the NCEP reanalysis. Although NWP reanalysis may have considerable uncertainties in simulating water cycle, they are often used to estimate moisture transport. Overall $\Theta_{SVR}$ correlates better with the NCEP reanalysis than $\Theta_{LR}$ and $\Theta_{NN}$ in seasonal variations. It also has the smallest RMS differences from the reanalysis data. The RMS bias of $\Theta_{LR}$ in the zonal

component is much larger than the other data sets. $\Theta_{LR}$ also shows large differences from the NCEP reanalysis for the zonal component over the Kuroshio and Gulf Stream extensions, and in the eastern Pacific and Atlantic south of the equator. $\Theta_{NN}$ has larger correlation discrepancies from the NCEP reanalysis in the meridional component. All three methods show poor correlation for the meridional components south of Baja California and west of Cape Verde of Africa, when the amplitude is weak and without clear annual cycles.

The advantage of SVR method in the retrieval of $\Theta$ from spacebased observations is clearly demonstrated. The rationale of the better performance by SVMs is discussed in many studies (e.g., Cortes & Vapnik, 1995; Gretton et al., 2001). It is briefly summarized here. First, the approach is relatively easy to use, because there are only a few parameters to adjust. The sparse setting of SVR, with the data training only based on the support vectors, avoids overfitting of the training data. By using the standard quadratic programming algorithms (Vapnik, 1995), only one global optimum is achieved. Mapping inputs into the high-dimensional feature space and introducing kernel function can solve the nonlinear relationship between inputs and outputs by turning a nonlinear regression to a linear fitting.

The present algorithm includes $u_{850}$ as an additional input parameter, which marginally improves the water balance over South America, compared with an earlier algorithm using SVR. Because $W$ can be measured by spacebased microwave radiometers, retrieval of $\Theta$ essentially becomes the problem of estimating $\mathbf{u}_e$. Directly mapping $\mathbf{u}_s$ to $\mathbf{u}_e$ only uses the vertically coherent part of wind profile. It is not sufficient over those areas when the surface wind is decoupled from the upper level and adding winds above the boundary layer should improve the methodology. The impact of adding the upper level

wind will be evaluated comprehensively over the global ocean, and optimization of the SVR algorithm will be further investigated.

## Acknowledgment

## References

Blanz, V., Schölkopf, B., Bülthoff, H., Burges, C., Vapnik, V., & Vetter, T. (1996). Comparison of view-based object recognition algorithms using realistic 3D models. *Artificial neural networks ICANN'96* (pp. 251−256).

Boser, B., Guyon, I., & Vapnik, V. N. (1992). A training algorithm for optimal margin classifiers. *Proceedings of 5th annual workshop on computer learning theory, Pittsburgh, PA* (pp. 144−152).

Bryan, F., & Oort, A. (1984). Seasonal variation of the global water balance based on aerological data. *Journal of Geophysical Research*, *89*, 11,717−11,730.

Burges, C. J. C. (1998). A tutorial on support vector machines for pattern recognition. *Data Mining and Knowledge Discovery*, *2*, 121−167.

Cherkassky, V., & Mulier, F. (1998). *Learning from data.* New York: John Wiley and Sons.

Cohen, J. L., Salstein, D. A., & Rosen, R. D. (2000). Interannual variability in the meridional transport of water vapor. *Journal of Hydrometeorology*, *1*, 547−553.

Cortes, C., & Vapnik, V. N. (1995). Support vector networks. *Machine Learning*, *20*, 273−297.

Drucker, H., Burges, C. J. C., Kaufman, L., Smola, A., & Vapnik, V. (1997). Support vector regression machines. *Advances in Neural Information Processing Systems*, *9*, 155−161.

Gretton, A., Doucet, A., Herbrich, R., Rayner, P. J. W., & Schölkopf, B. (2001). Support vector regression for black-box system identification. *Statistical signal processing, proceedings of the 11th IEEE signal processing workshop* (pp. 341−344).

Hardle, W. (1990). *Applied nonparametric regression: Cambridge University Press 333 pp.*

Hayden, C. M., & Pursor, R. J. (1995). Recursive filter objective analysis of meteorological fields: Applications to NESDIS operational processing. *Journal of Applied Meteorology*, *34*, 3−15.

Heta, Y., & Mitsuta, Y. (1993). An evaluation of evaporation over the tropical Pacific Ocean as observed from satellites. *Journal of Applied Meteorology*, *32*, 1242−1247.

Horvath, A., & Davies, R. (2001). Feasibility and error analysis of cloud motion wind extraction from near-simultaneous MISR measurements. *Journal of Atmospheric and Oceanic Technology*, *18*, 591−680.

Joachims, T. (1999). Making large-scale SVR learning practical. In B. Schölkopf, C. Burges, & A. Smola (Eds.), *Advances in kernel methods— Support vector learning* : MIT Press.

Liu, W. T. (1987). *1982–1983 El Nino atlas—Nimbus-7 microwave radiometer data.JPL Publication*, vol. 87-5. (pp. ): Jet Propulsion Laboratory 68 pp.

Liu, W. T. (1993). Ocean surface evaporation. In R. J. Gurney, J. Foster, & C. Parkinson (Eds.), *Atlas of satellite observations related to global change* (pp. 265−278). Cambridge: Cambridge University Press.

Liu, W. T. (2002). Progress in scatterometer application. *Journal of Oceanography*, *58*, 121−136.

Liu, W. T., & Tang, W. (2005). Estimating moisture transport over ocean using spacebased observations from space. *Journal of Geophysical Research*, *110*, D10101. doi:10.1029/ 2004JD005300

Liu, W. T., Tang, W., & Niiler, P. P. (1991). Humidity profiles over ocean. *Journal of Climate*, *4*, 1023−1034.

Liu, W. T., Xie, X., Tang, W., & Zlotnicki, V. (2006). Spacebased observations of oceanic influence on the annual variation of South American water balance. *Geophysical Research Letters*, *33*, L08710. doi:10.1029/2006GL025683

Meyer, D., Leisch, F., & Hornik, K. (2003). The support vector machine under test. *Neurocomputing*, *55*, 169−186.

Mo, K. -C., & Higgins, R. W. (1996). Large-scale atmospheric moisture transport as evaluated in the NCEP/NCAR and the NASA/DAO reanalyses. *Journal of Climate*, *9*, 1531−1545.

Müller, K. -R., Smola, A., Rätsch, G., Schölkopf, B., Kohlmorgen, J., & Vapnik, V. (1997). Predicting time series with support vector machines. *Proceedings of international conference on artificial neural networks* (pp. 999).

Peixoto, J. P., Salstein, D. A., & Rosen, R. D. (1981). Intra-annual variations in large-scale moisture fields. *Journal of Geophysical Research*, *86*, 1255−1264.

Rasmusson, E. M. (1967). Atmospheric water vapor transport and the water balance of North America. *Monthly Weather Review*, *95*, 403−426.

Roads, J. O., Chen, S. -C., Kao, J., Langley, D., & Glatzmaier, G. (1992). Global aspects of the Los Alamos general circulation model hydrological cycle. *Journal of Geophysical Research*, *97*, 10,051−10,068.

Rosen, R. D., Salstein, D. A., & Peixoto, J. P. (1979). Variability in the annual fields of large-scale atmospheric water vapor transport. *Monthly Weather Review*, *107*, 26−37.

Schölkopf, B. (1997). *Support vector learning.* München: R. Oldenbourg Verlag Doktorarbeit, TU Berlin. http://www.kernel-machines.org/

Smola, A. J., & Schölkopf, B. (2004). A tutorial on support vector regression. *Statistics and Computing*, *14*, 199−222.

Starr, V. P., & White, R. M. (1955). Direct measurement of the hemispheric poleward flux of water vapor. *Journal of Marine Research*, *14*(3), 217−225.

Tang, B., & Mazzoni, D. (2006). Multiclass reduced-set support vector machines. *Proceedings of the 23th international conference on machine learning (ICML)*.

Trenberth, K. E., & Guillemot, C. J. (1995). Evaluation of the global atmospheric moisture budget as seen from analyses. *Journal of Climate*, *8*, 2255−2272.

Vapnik, V. N. (1995). *The nature of statistical learning theory: Springer.*

Vapnik, V. N. (1998). *Statistical learning theory.* New York: Wiley.

Wentz, F. J. (1997). A well-calibrated ocean algorithm for special sensor microwaver/imager. *Journal of Geophysical Research*, *102*(C4), 8703−8718.

Xie, X., & Liu, W. T. (2005). Hydrological budget in the Tropical Pacific. *16th conference on climate variability and change* : Amer. Meteor. Soc. http://ams.confex.com/ams/pdfpapers/85220.pdf