

Optimizing Machine Learning Algorithms for Hyperspectral Very Shallow Water (VSW) Products

W. Paul Bissett

Florida Environmental Research Institute

10500 University Center Dr.

Suite 140

Tampa, FL 33612 USA

phone: (813) 866-3374 x102 fax: (813) 977-8057 email: pbissett@feriweb.org

Award Number: N000140810622

<http://www.FERIweb.org>

http://www.onr.navy.mil/sci_tech/32/322/ocean_optics_biology.asp

LONG-TERM GOALS

This one-year effort will focus on the transition of FERI's machine learning algorithms for HyperSpectral Imagery (HSI) in the VSW into a distributable code set. This will provide a stable code platform for the application and transition of machine learning-based hyperspectral classification techniques into 6.3/6.4 programs. (This work was funded mid-year 2008.)

OBJECTIVES

Our objective is to focus on three areas of application research and transitions. First, we will transition our machine learning-based algorithms and computer code for the determination of bathymetry, bottom type, and water column Inherent Optical Properties from HyperSpectral Imagery (HSI) into a deliverable Message Passing Interface (MPI) program that may be easily used by other research and military operators. Second, we will use this program to determine the impacts of the granularity of the classification database on the inversion bathymetry, bottom type, and IOPs. Third, we will move beyond the use of single pixel HSI inversion to the use of spatial context-filtering to remove pixel-to-pixel noise inherent in the HSI data.

APPROACH

Task 1

In previous works, a Look-Up Table (LUT) algorithm was used in accurately predicting bathymetry (Mobley et al. 2002, Bissett et al. 2004, Bissett et al. 2005, Mobley et al. 2005, Lesser and Mobley, 2008). The LUT approach is a subset of a larger body of artificial intelligence work concerned with algorithms and techniques that "teach" machine to learn from the examination of data and rules. This body of work is aptly called "machine learning" and some of its techniques include decision trees, genetic algorithms, and neural networks. The LUT approach is a subset of the k-Nearest Neighbor (kNN) algorithm, which is in the family of supervised learning algorithms.

Our use of the kNN algorithm maps a single HSI remote sensing reflectance vector, $Rrs(\lambda)$, onto a database of estimated $Rrs(\lambda)$. This database is created by providing the attributes of bathymetry,

Report Documentation Page

Form Approved
OMB No. 0704-0188

Public reporting burden for the collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing this burden, to Washington Headquarters Services, Directorate for Information Operations and Reports, 1215 Jefferson Davis Highway, Suite 1204, Arlington VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to a penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

1. REPORT DATE 2008		2. REPORT TYPE		3. DATES COVERED 00-00-2008 to 00-00-2008	
4. TITLE AND SUBTITLE Optimizing Machine Learning Algorithms for Hyperspectral Very Shallow Water (VSW) Products				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES) Florida Environmental Research Institute, 10500 University Center Dr., Tampa, FL, 33612				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT Approved for public release; distribution unlimited					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES	19a. NAME OF RESPONSIBLE PERSON
a. REPORT unclassified	b. ABSTRACT unclassified	c. THIS PAGE unclassified			

spectral bottom reflectance, and spectral IOPs to the radiative transfer routines of Ecolight (which is a high speed variant of Hydrolight, Mobley, 1994). We select the classification of the measured R_{rs} vector based on the best match of measured $R_{rs}(\lambda)$ to estimated $R_{rs}(\lambda)$. The LUT algorithm is based on a single best fit for our classification, i.e. $k = 1$. However, more recent work suggested that we could achieve a better classification by selecting a larger number for k , e.g. $k = 50$ (Bissett et al. 2006a). This larger number for k provides better accuracy and precision, as well as provides us with the ability to create confidence intervals for our classifications of bathymetry.

When classifying new spectra, the distance or angle between each measured spectrum and estimated spectrum in the database is calculated. The k nearest neighbors to that spectra (those having the smallest distances or angles), are considered sufficiently qualified to predict the corresponding attributes of bathymetry, bottom type, and IOP set. We have used the following metrics for the calculation of distance (Euclidean, Manhattan, Chebyshev, Canberra and Bray Curtis) and/or angle (Angular Separation and Correlation Coefficient). In general, our applications suggest that the Manhattan distance and the Correlation Coefficient angle metrics to be the best metrics to use for this algorithm. Once the set of nearest neighbors are determined, the attribute (e.g. bathymetry) of a pixel may be determined by a majority vote from the k nearest neighbor vectors. In the event of a tie, a prediction is made randomly from amongst the majority classes.

The computer code used in our creation of the estimated $R_{rs}(\lambda)$ database and the spectral matching of the measured versus estimated $R_{rs}(\lambda)$ is functional for scientific research; however it not well developed for transition for use by others in testing and evaluation applications. **Our first task of this project will build upon our past research efforts to provide a Message Passing Interface (MPI) executable version of our kNN workbench for the inversion of hyperspectral imagery.** This code will be distributed to research and military partners for testing and evaluation purposes, as well as to complete Task 2 and 3.

Task 2

The spectrum for one particular depth, bottom type, and set of inherent optical properties may closely match a multitude of spectra with many different attributes (Figure 1). The selection of a single nearest neighbor may produce noisy predictions because of the noise in both the measured and estimated $R_{rs}(\lambda)$. The total prediction noise is a function of the noise associated with the measured $R_{rs}(\lambda)$, which contains components of sensor and environmental noise, and the noise associated with the estimation of $R_{rs}(\lambda)$ in the training database. This noise is evident in the “speckling” that may be associated with these inversion techniques (Figure 2). The use of kNN algorithms work to reduce noise of the prediction by increasing the probability that a spectrum presented for classification will come from the majority class of proximally-located spectral vectors, rather than a single “lucky” spectrum. In this case, rather than selecting the single database spectrum “O” that is closest to the measured spectrum (represented by the square in Figure 1), a majority vote of all of the nearest neighbors around the square is used to make the prediction of the attribute (e.g. bathymetry) at that pixel location. Choosing the majority class creates a less variable space from which to make a decision, making it is less likely to produce different classifications due to small amounts of noise in the spectra.

However, as the size of the training database increases (through the increase in number of bathymetry depths, bottom types, or IOP sets) the number of nearest neighbors also increases (Figure 3). This in turn causes a problem with “non-uniqueness” in the selection of the appropriate class, and its component attribute. This, in turn, causes increasing noise in the map of the estimated attribute (e.g.

bathymetry), and therefore it became very important to have the appropriate “granularity”, or the proper step size in the discrete selection of attributes that are used in the creation of the training database. In this specific case, it means that we need to be selective in the selection of number of depth levels, bottom types, and IOP sets that we use to create the estimated $Rrs(\lambda)$ database. **The second Task of this project will be to use the code from Task 1 to rapidly test the impacts of granularity of attribute selection on the accuracy and precision of bathymetry estimated from our kNN code and the HSI data from Horseshoe Reef and St. Joseph Bay, FL¹ (Bissett et al. 2006b).**

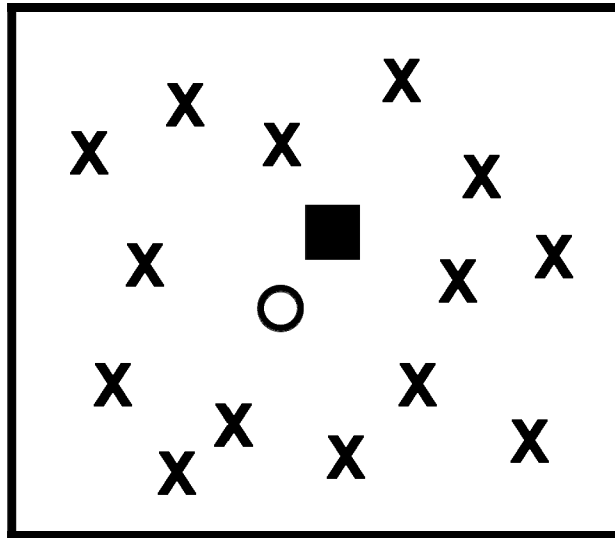


Figure 1. Xs and Os are the classes of examples belonging to the training database. The measured spectrum, λ , is closer to the O than any X. In kNN, multiple nearest neighbors are used to vote on the appropriate class. If $k = 1$, class O is chosen. If $k > 1$, a vote amongst all the classes X is chosen. The total number of X is dependent on the value of k , and in which will include O in the retrieved set. The estimate of the attribute may then be calculated from any number of statistical calculations on the set of Xs, e.g. mean, majority vote, etc.

Task 3

The problem of sensor and environmental noise is a critical issue in the retrieval of accurate bathymetry from maps of HSI data. There are many sources of environmental noise in the collection of sensor measured radiance, for example surface waves that alter the reflection surface and path length to the bottom reflectance target. These surface noise effects are commingled with the atmospheric and illumination correction noise to produce spatially varying $Rrs(\lambda)$ over areas with identical bathymetry, bottom types, and IOPs (Figure 2). In order to reduce the impacts of this environmentally generated noise component, we should use the spatial context of the measured spectrum during the selection of the nearest neighbor classes, and subsequent estimate of the attribute of interest.

¹ The use of St. Joseph Bay, FL data will depend on acquiring accurate bathymetry from the State of Florida. If we do not receive bathymetry of sufficient quality, we will focus on the Horseshoe Reef imagery.

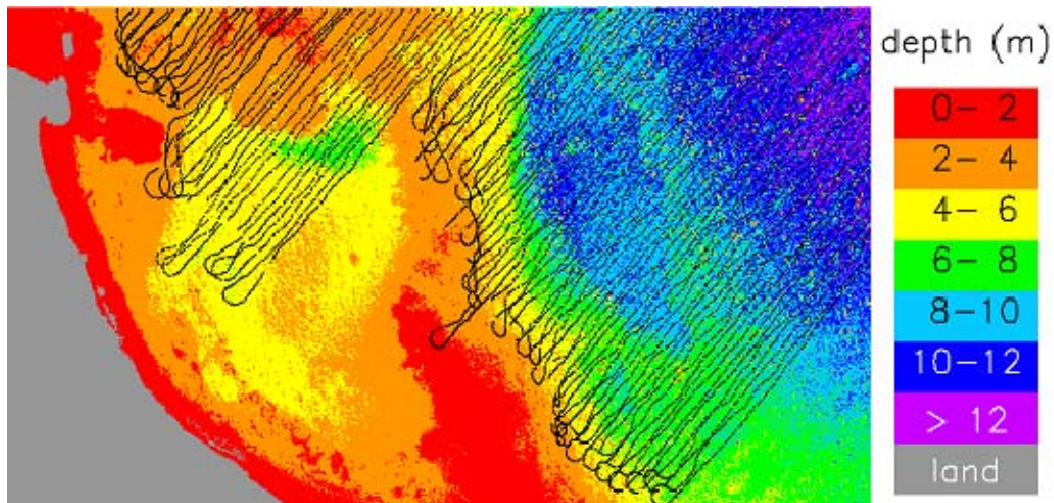


Figure 2. LUT bathymetry estimate for Horseshoe Reef, Bahamas. The black dots show the locations of the acoustic pings. The color-coded depths are for the unconstrained LUT retrieval ($k = 1$) applied to the entire image. The speckling in bathymetry is evident throughout the image.

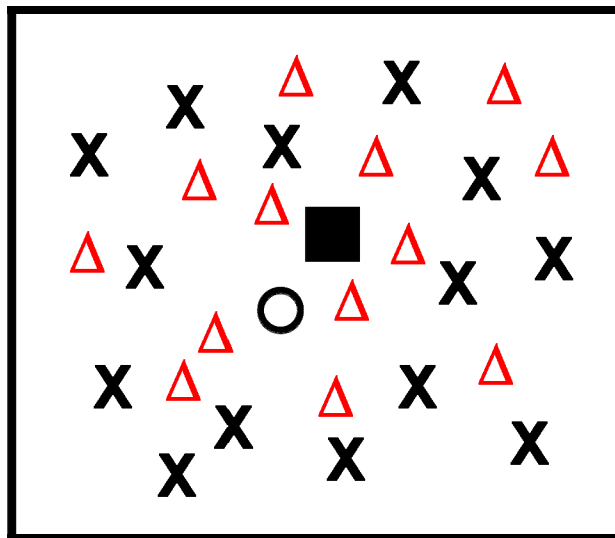


Figure 3. Xs and Os are the classes of examples belonging to the training database and are the same as Figure 1. The Δ 's are addition classes resulting from increasing the depth resolution, as well as the number of bottom types and IOP sets. In this case discussed in the text, these Δ 's may contain attributes that are unrepresentative of the actual values and represents a non-unique solution to this inversion problem. The selection of the appropriate depth intervals or range of bottom types and IOPs sets is important to reducing this non-uniqueness. The term granularity is used to describe the separation between the discrete levels in the attributes.

Heretofore we have done point- or pixel-specific classification of HSI data. That is, each pixel is classified (for depth, bottom type, and water IOPs) independently of its neighbors, and only the spectral character of the pixel is used in its classification. Task 3 will be to evaluate spatial context-sensitive classification, which means that we will incorporate information about the spatial neighborhood (the spatial context) of a pixel to assist with its classification. Context-sensitive classification is often used in traditional terrestrial thematic mapping (e.g., Richards and Jia, 2006, §8.8) and some of those techniques may be beneficial for our oceanic problem.

This Task will evaluate two types of context-filtering – (1) pre-filtering of the $Rrs(\lambda)$ spectra before classification, and (2) context-filtering of the retrieved attributes after classification. The first type of context-filtering seeks to reduce the noise in $Rrs(\lambda)$ spectra by replacing the spectrum value at each wavelength with the median value of the spectra in a spatial area surrounding the pixel of interest, say a 3 x 3 grid of pixels centered on the one of interest. This spatial filter is applied wavelength by wavelength. At wavelengths where $Rrs(\lambda)$ is mostly signal, the final spectrum will not change by much. At wavelengths where $Rrs(\lambda)$ is noisy, the noise in the surrounding pixels will tend to average out and the final spectrum values over the entire image area will be less noisy than the original.

The second type of context-filtering involves post-processing the retrievals themselves, rather than the original image spectra. In the case of real numbered attributes, such as bathymetry, we can apply a median filter to the retrieved depth. For bottom type and IOP set, the way forward is less clear. Each of these attributes is assigned a type with a specific vector (or set of vectors in the case of IOPs) of spectral values. How we filter “Dark Sediment” with “Sparse Vegetation” or “Highly absorbing and scattering waters #1” with “Case 1, chlorophyll a = 0.5 mg m⁻³” will be a challenge. It may require some iterative solution that context-filters bathymetry first, and solves the kNN again using a constrained bathymetry solution approach. It may also be highly dependent on the granularity study in Task 2. These are the issues that we will address in this Task.

WORK COMPLETED

Task (1) has been completed and the serial and MPI versions of our optimized machine learning code is available for v 0.1.0 release. The code will be distributed in a generic Red Hat Package Manager (RPM; http://en.wikipedia.org/wiki/RPM_Package_Manager) format for installation on Red Hat, Fedora, and CentOS version of Linux.

IMPACT/APPLICATIONS

This effort will deliver an application for testing and evaluation of our machine learning approaches to bathymetry estimation in Very Shallow Waters (VSW). While it is being demonstrated on hyperspectral imagery, the techniques and computer code may be used with any set of spectral reflectance data. As such the deliverables from this effort will allow other to create maps of depths, bottom types, and water clarity from a variety of airborne and space-based spectral sensors planned for operational deployment.

RELATED PROJECTS

This work is being conducted in conjunction with Dr. Curtis D. Mobley at Sequoia Scientific, Inc., who is funded under this effort for the collaboration. These techniques developed here are now being

applied to imagery of Australian coastal waters in a comparison of several different hyperspectral remote sensing algorithms for a variety of environments. That comparison study is being led by A. Dekker of CSIRO.

REFERENCES

Bissett, W.P., DeBra, S., Kadiwala, M., Kohler, D., Mobley, C., Steward, R., Weidemann, A., Davis, C.O., Lillycrop, J. and Pope, R., 2004. Development, validation, and fusion of high resolution active and passive optical imagery. Ocean Optics XVII, Fremantle, AU.

Bissett, W.P., DeBra, S., Kadiwala, M., Kohler, D.D.R., Mobley, C.D., Steward, R.G., Weidemann, A.D., Davis, C.O., Lillycrop, J. and Pope, R.L., 2005. Development, validation, and fusion of high-resolution active and passive optical imagery. In: I. Kadar (Editor), Signal Processing, Sensor Fusion, and Target Recognition XIV. Proceedings of SPIE Vol. 5809. SPIE, Bellingham, WA, pp. 341-349.

Bissett, W.P., Banfield, R., Kohler, D.D.R. and Mobley, C.D., 2006a. Ascribing Confidence Intervals to HyperSpectral Imaging (HSI) Bathymetry Maps, Ocean Optics XVIII, Montreal, Quebec, Canada.

Bissett, W.P. and Kohler, D.D.R., 2006. St. Joseph Bay Aquatic Preserve Hyperspectral Imaging - FINAL REPORT, Florida Environmental Research Institute, Tampa, FL. Florida Department of Environmental Protection, Contract Number RM055 (<http://www.feriweb.org/pubs/index.html#tech>)

Lesser, M. P. and C. D. Mobley. 2007. Bathymetry, optical properties, and benthic classification of coral reefs using hyperspectral remote sensing imagery. *Coral Reefs*, (26) 819-829.

Mobley, C.D., 1994. Light and Water. Academic Press, San Diego, CA, 592 pp.

Mobley, C.D., Sundman, L., Davis, C.O., Montes, M. and Bissett, W.P., 2002. A look-up-table approach to inverting remotely sensed ocean color data, Ocean Optics XVI. Office of Naval Research Ocean, Atmosphere, and Space S&T Department, Santa Fe, NM.

Mobley, C. D., L. K. Sundman, C. O. Davis, T. V. Downes, R. A. Leathers, M. J. Montes, J. H. Bowles, W. P. Bissett, D. D. R. Kohler, R. P. Reid, E. M. Louchard, and A. Gleason, 2005. Interpretation of hyperspectral remote-sensing imagery via spectrum matching and look-up tables, *Appl. Optics*, 44(17) 3576-3592.

Richards, J. A. and X. Jia, 2006. *Remote Sensing Digital Image Analysis, 4th Edition*. Springer, 439 pages.