

Spatially Explicit Predictors of Indicators of Water Quality: Example from Wadeable Streams in the U.S.

Mostafa Shirazi and Marc Weber
Western Ecology Division
U.S. Environmental Protection Agency
Corvallis, Oregon

Acknowledgement

This research is derived from four national-scale databases which were produced over decades by hundreds of staffs of government agencies and their collaborators from the academic and research institutions, efforts which inspire a heartfelt gratitude.

- NRCS National Soil Database SSURGO and STATSGO
- USGS National Land Cover Database (NLCD)
- USEPA National Aquatic Resource Database, Wadeable Stream Assessment (WSA).
- USEPA Ecoregions of the United States

Objective (1)

Spatially Explicit Prediction of Water Quality

- The Environmental Protection Agency (USEPA), in collaboration with the States, is required to assess and report the condition of surface waters nation-wide.
- The Wadeable Stream Assessment (WSA), is a nation-wide probability-based survey, which located a number of wadeable stream sites, applied uniform field procedures to each site, and collected data, which developed water quality indicators (WQI).
- Statistically valid inferences about all surface waters can be drawn from the sampled sites, but an explicit connection of WQI with the ecosystem and its land use, which influence water quality is missing.
- This study uses soil and land use predictors to establish explicit linkage between *WQI* of WSA watersheds and the Ecological Regions of the conterminous United States.

Objective (2)

Using SC+LC in predictive models

- Describe SSURGO and STATSGO soil *map units* by the mean Soil Characteristics (SC) of soil layers and soil components.
- Intersect soil map units with NLCD (LC) to produce a combined national SC+LC coverage: *SL map units*.
- Intersect *SL map units* with the following, and describe them identically by their mean SC+LC:
 - 1- 1392 WSA “Probability” watersheds,
 - 2- 451 “least Impacted, Reference” watersheds, and
 - 3- 967 US Ecoregions.
- Use SC+LC in models as *independent predictors* of WQI of a *not-tested stream*, while using the remaining WQI as *co-predictors*.

Selected SC and LC of Soil Map Units and Water Quality indicators (WQI) of WSA Sites

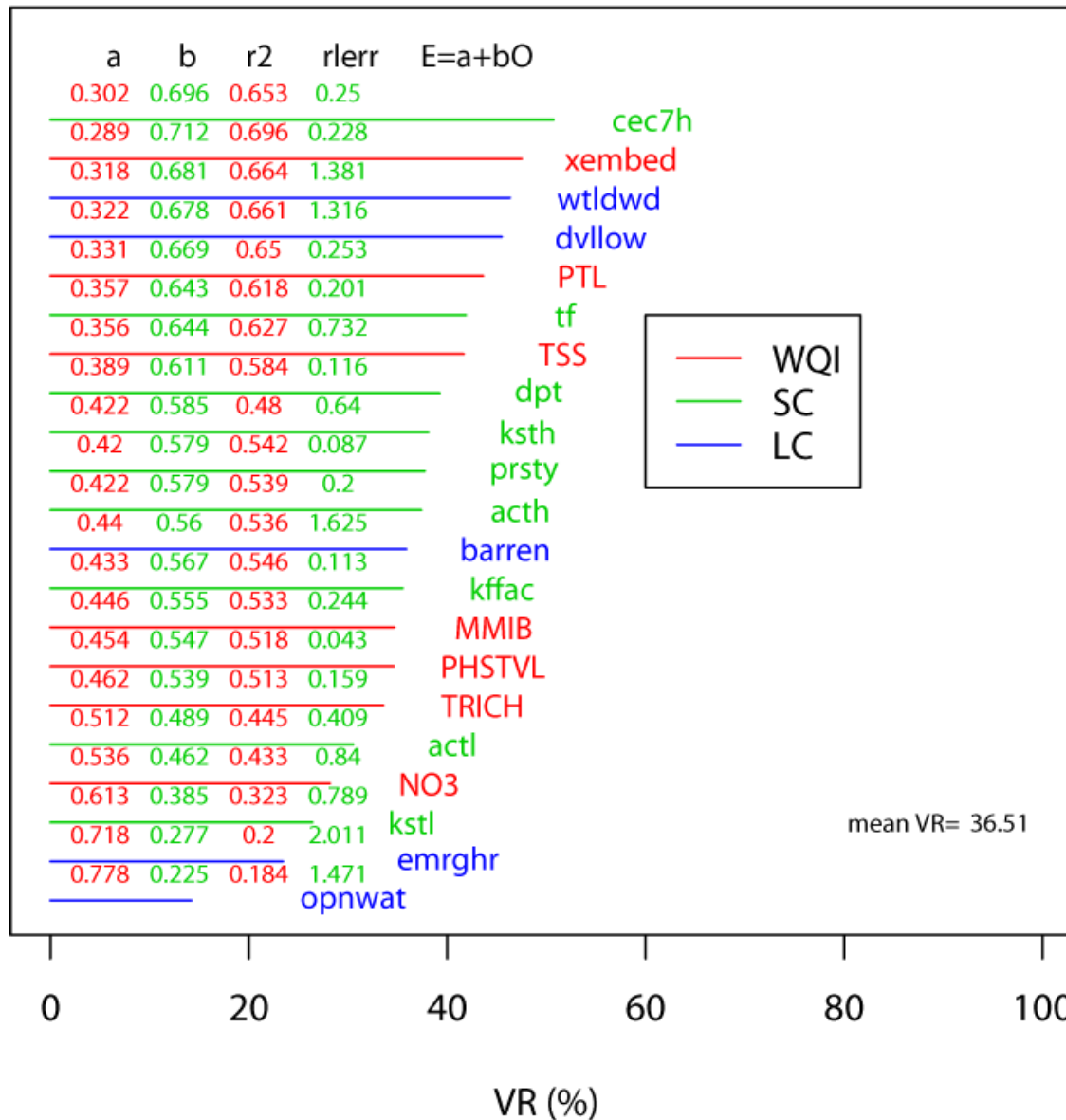
- **Soil Characteristics SC** (57 SC from 281) Depth, Texture, Bulk Density, Available Water Capacity, Organic Matter, pH, CEC,...., Elevation, Slope, Air Temperature, Precipitation, Frost Free Days,...
- **Land Cover LC** (12 classes from 29), developed space, cropland, pastureland, deciduous forests, evergreen forests,...
- **Stream Water Quality Indicators WQI** (17 items selected): pH, ANC, DOC, Stream Substrate, Taxa Richness, PTL, NO₃, SO₄, CL.

Models and Test Strategies (1)

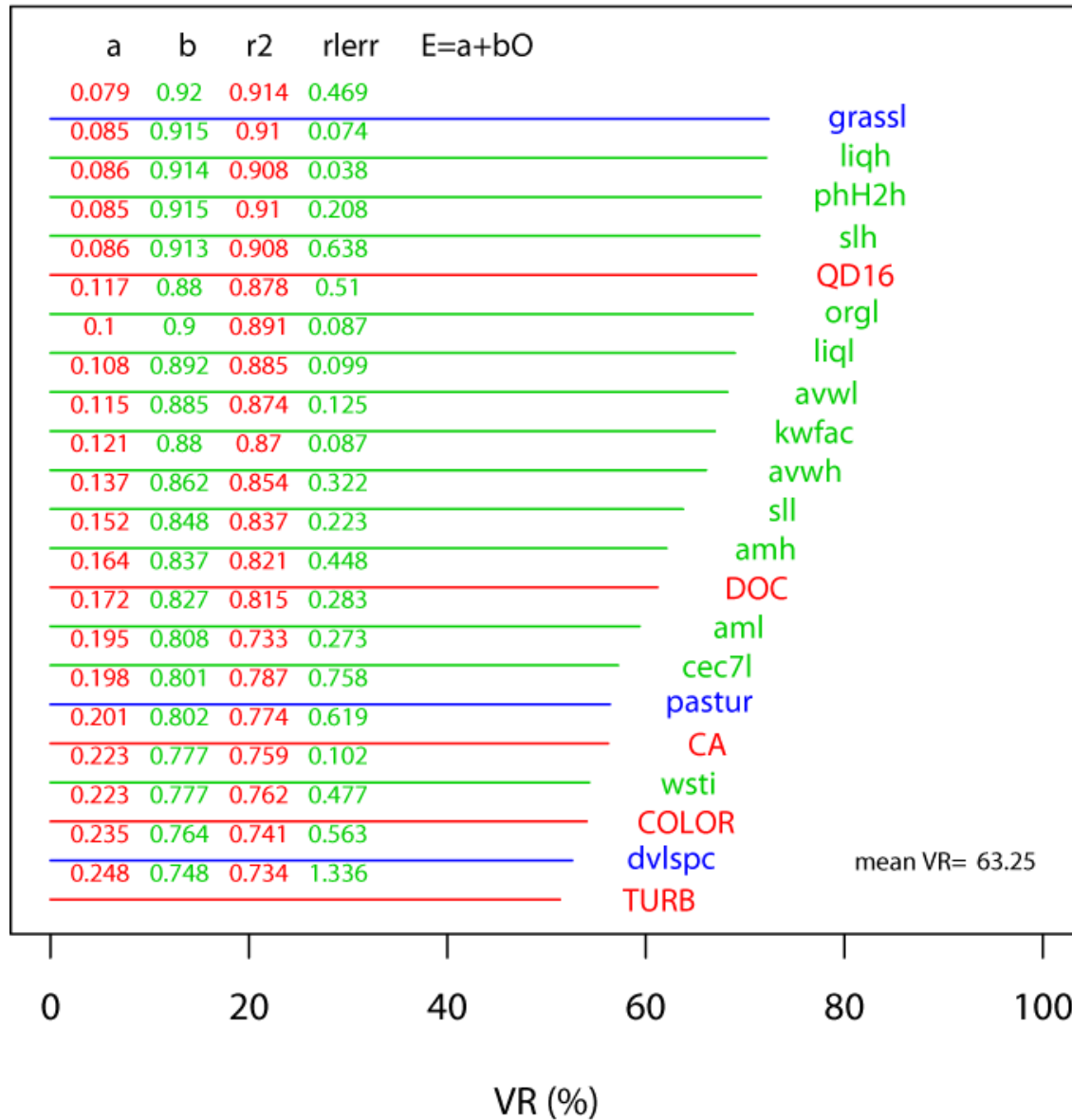
Rank and Predictability of SC, LC, and WQI

- Compare the sample (V_s) with the conditional variance (V_c), of each vector in the list SC, LC, and WQI, and the % variance reduction (VR) achieved because of vector's spatial relationships with the others not used from the list: **$VR = 100[1.0 - (V_c / V_s)^{1/2}]$** .
- Use one of the 1392 watersheds as a *set-aside*, build a model using the remaining 1391, and predict, in an extrapolation mode, each one of the *set-aside's* SC, LC, and WQI. Repeat 1392 times to calculate *relative standard error*, and r^2 .

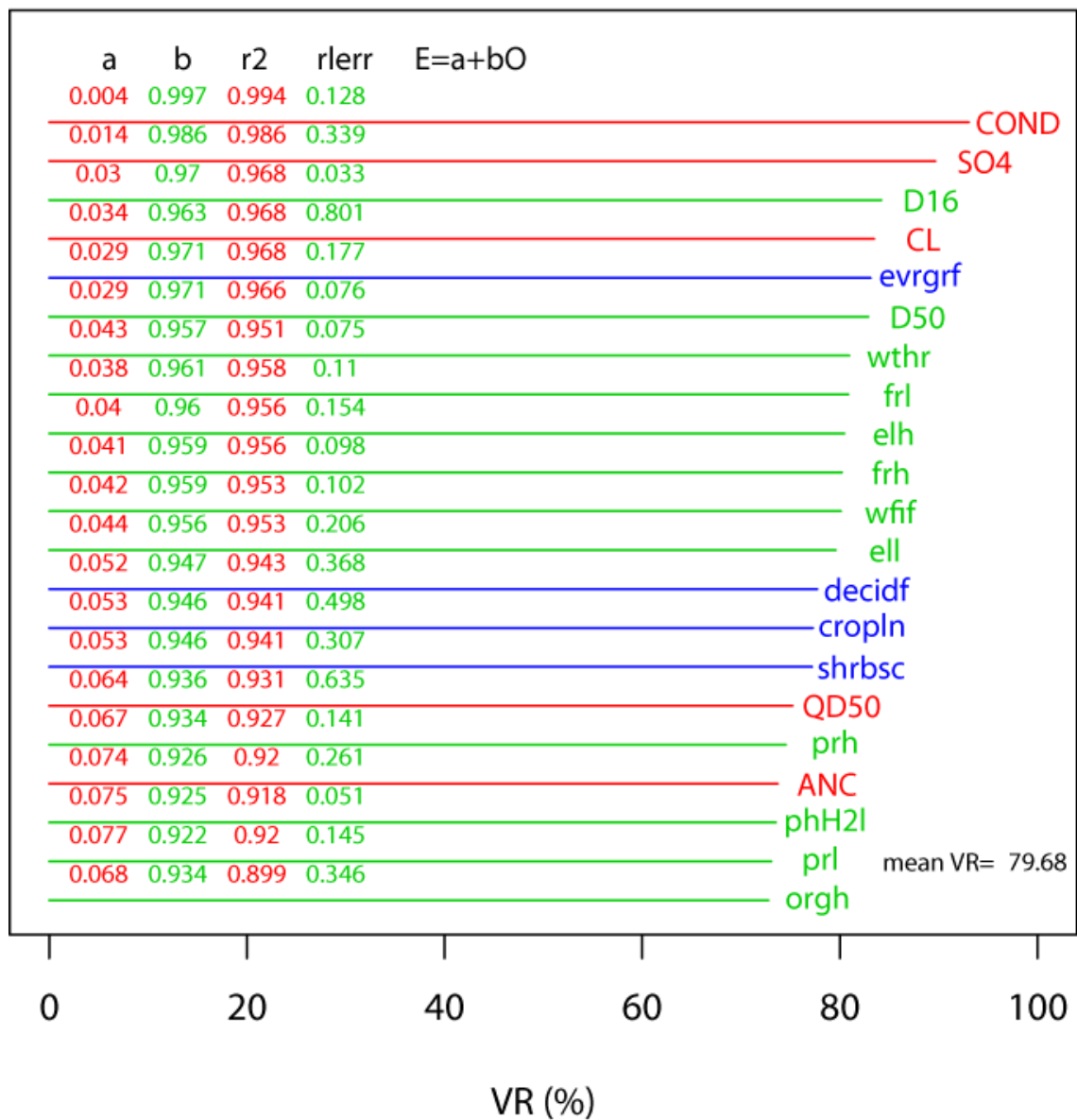
WQI, SC, LC



WQI, SC, LC



WQI, SC, LC

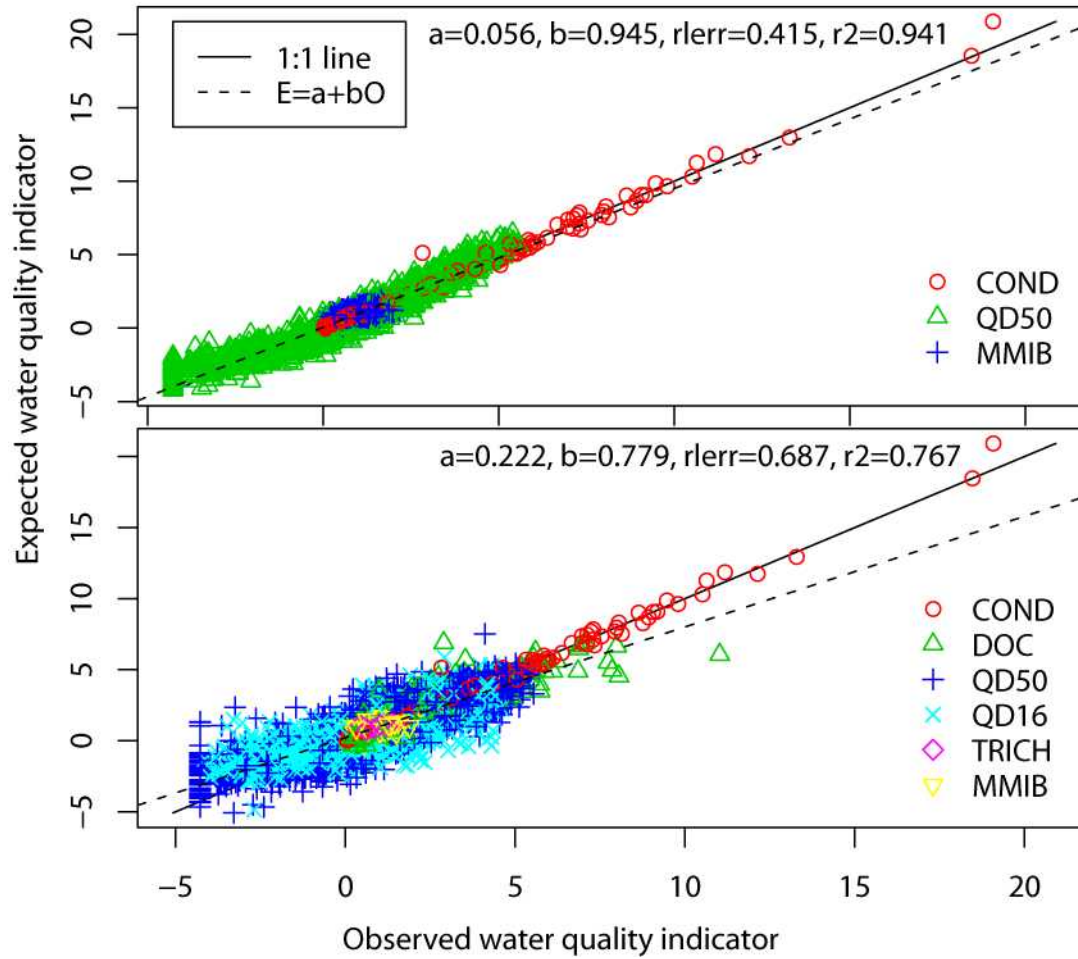


Models and Test Strategies (2)

Joint Prediction of WQI (small and large VR)

- A weakly predicted WQI (mean VR=37%) improves in a joint prediction with one or more strongly predicted WQI (mean VR=80%).
- An example of joint prediction tested with data from 1392 set-aside watershed sites (next slide).

Jointly predicted WQI tested with 1392 set-aside watershed data

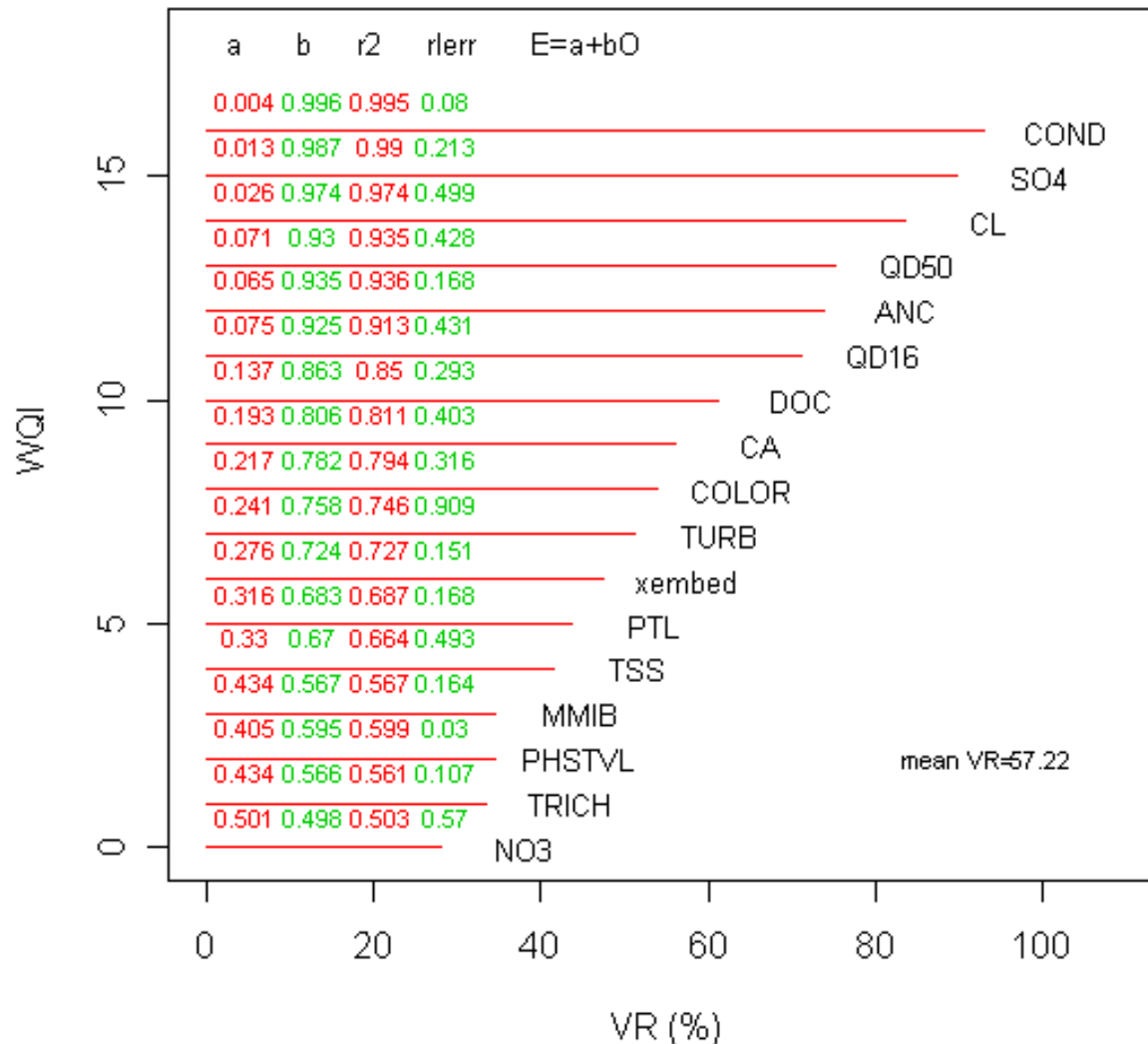


Models and Test Strategies (3)

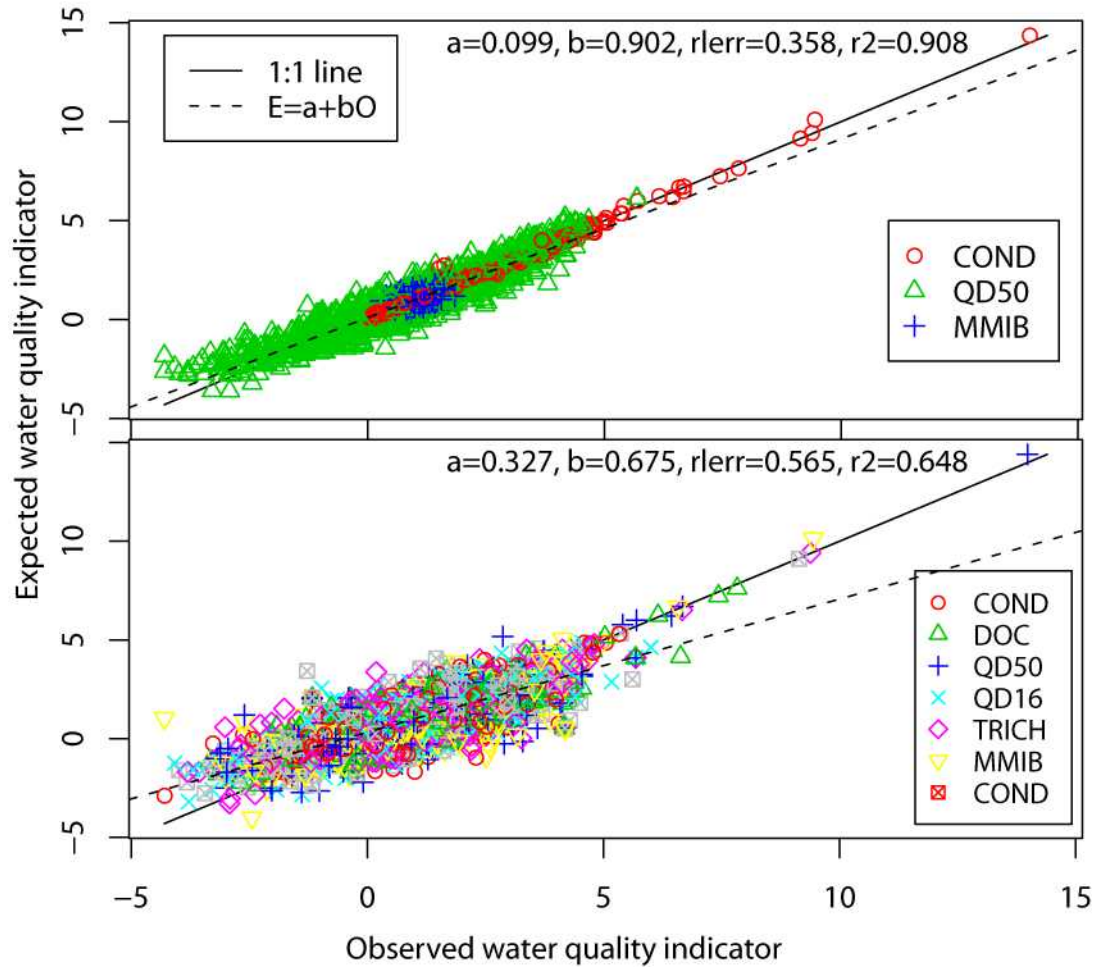
Using Neighbors of a Set-aside Watershed

- **1-** Test models using 1392 probability sites with *WQI* obtained indirectly, from the “neighbor” watersheds of a set-aside “host” watershed.
- **2-** Locate the nearest neighbors of each host watershed using a sorted list of 1392 *squared generalized distance* (*SL_sqrd*) of *SC+LC*.
- **3-** Predict *WQI* singly and jointly

Predicted WQI using nearest neighbor



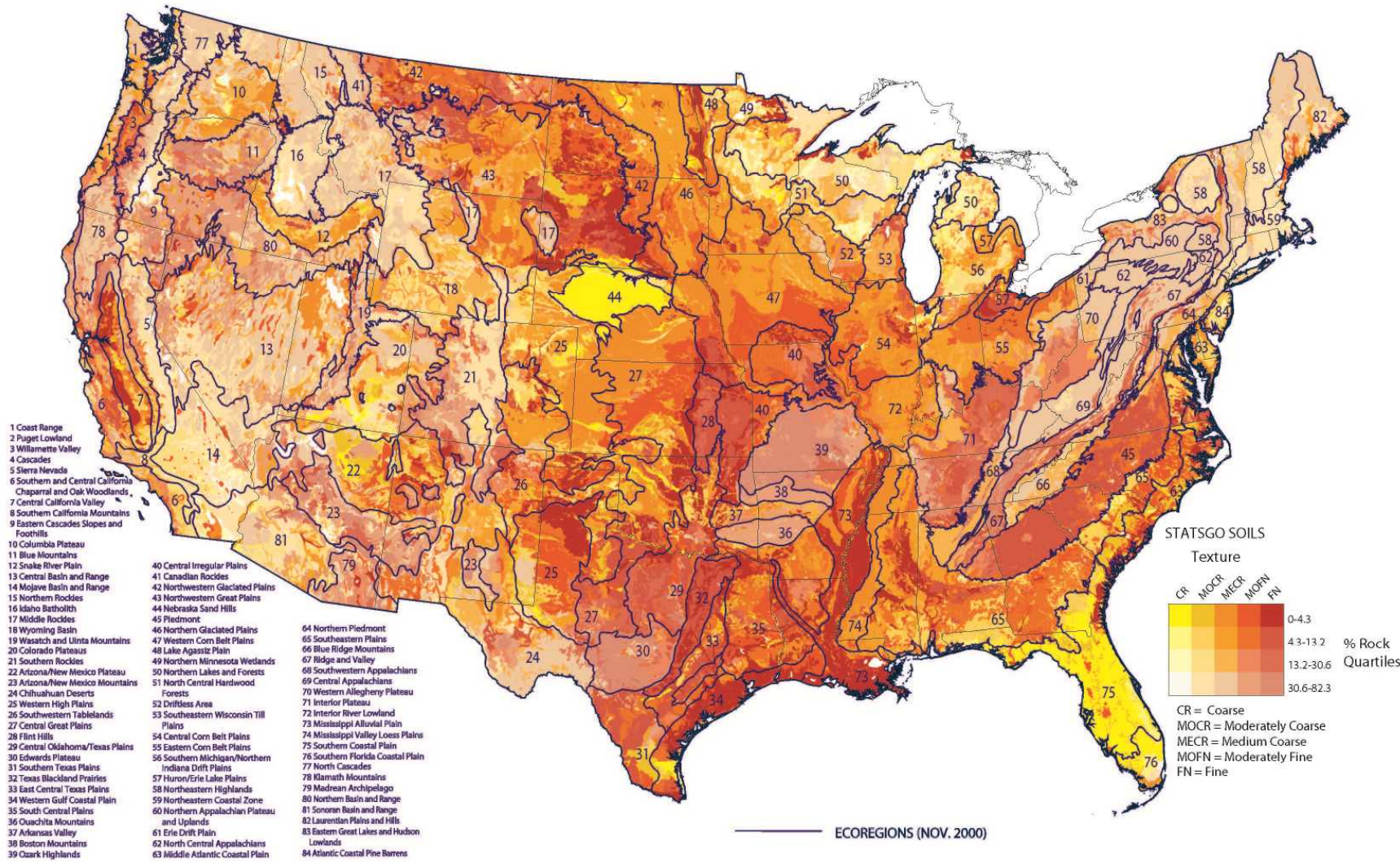
"Neighbors" of 1392 set-aside watersheds donate WQI to "host"



Linking WSA Watersheds and Ecoregions (1)

Example of quantification: Using 26 Soil Characteristics (SC)

- Ecological regions consist of ecosystems (forests, wetlands, deserts) which evolve under regional climate and geology, maintain particular water quality, sustain plant and animal communities, and variously invite human land use.
- Extensive mapped data (including land use), professional judgment, and interpretations determined the ecological boundaries in the U.S. for level *III* and level *IV* classifications, which are available online: (<http://www.epa.gov/wed/pages/ecoregions.htm>)
- A complete quantitative description of ecoregions is not feasible, but previous quantification, using, as a first approximation, 26 soil characteristics from STATSGO, variously matched boundaries of Level III ecoregions, for example, the soil texture map (next slide).



- 1 Coast Range
- 2 Puget Lowland
- 3 Willamette Valley
- 4 Cascades
- 5 Sierra Nevada
- 6 Southern and Central California Chaparral and Oak Woodlands
- 7 Central California Valley
- 8 Southern California Mountains
- 9 Eastern Cascades Slopes and Foothills
- 10 Columbia Plateau
- 11 Blue Mountains
- 12 Snake River Plain
- 13 Central Basin and Range
- 14 Mojave Basin and Range
- 15 Northern Rockies
- 16 Idaho Batholith
- 17 Middle Rockies
- 18 Wyoming Basin
- 19 Wasatch and Uinta Mountains
- 20 Colorado Plateau
- 21 Southern Rockies
- 22 Arizona/New Mexico Plateau
- 23 Arizona/New Mexico Mountains
- 24 Chihuahuan Deserts
- 25 Western High Plains
- 26 Southwestern Tablelands
- 27 Central Great Plains
- 28 Flint Hills
- 29 Central Oklahoma/Texas Plains
- 30 Edwards Plateau
- 31 Southern Texas Plains
- 32 Texas Blackland Prairies
- 33 East Central Texas Plains
- 34 Western Gulf Coastal Plain
- 35 South Central Plains
- 36 Ouachita Mountains
- 37 Arkansas Valley
- 38 Boston Mountains
- 39 Ozark Highlands
- 40 Central Irregular Plains
- 41 Canadian Rockies
- 42 Northwestern Glaciated Plains
- 43 Northwestern Great Plains
- 44 Nebraska Sand Hills
- 45 Piedmont
- 46 Northern Glaciated Plains
- 47 Western Corn Belt Plains
- 48 Lake Agassiz Plain
- 49 Northern Minnesota Wetlands
- 50 Northern Lakes and Forests
- 51 North Central Hardwood Forests
- 52 Driftless Area
- 53 Southeastern Wisconsin Till Plains
- 54 Central Corn Belt Plains
- 55 Eastern Corn Belt Plains
- 56 Southern Michigan/Northern Indiana Drift Plains
- 57 Huron/Erie Lake Plains
- 58 Northeast Highlands
- 59 Northeastern Coastal Zone
- 60 Northern Appalachian Plateau and Uplands
- 61 Erie Drift Plain
- 62 North Central Appalachians
- 63 Middle Atlantic Coastal Plain
- 64 Northern Piedmont
- 65 Southeastern Plains
- 66 Blue Ridge Mountains
- 67 Ridge and Valley
- 68 Southwestern Appalachians
- 69 Central Appalachians
- 70 Western Allegheny Plateau
- 71 Interior Plateau
- 72 Interior River Lowland
- 73 Mississippi Alluvial Plain
- 74 Mississippi Valley Loess Plains
- 75 Southern Coastal Plain
- 76 Southern Florida Coastal Plain
- 77 North Cascades
- 78 Klamath Mountains
- 79 Madream Archipelago
- 80 Northern Basin and Range
- 81 Sonoran Basin and Range
- 82 Laurentian Plains and Hills
- 83 Eastern Great Lakes and Hudson Lowlands
- 84 Atlantic Coastal Pine Barrens

STATSGO SOILS

Texture

CR	MOCR	MECR	MOFN	FN
0-4.3				
4.3-13.2				
13.2-30.6				
30.6-82.3				

% Rock
Quartiles

CR = Coarse
MOCR = Moderately Coarse
MECR = Medium Coarse
MOFN = Moderately Fine
FN = Fine

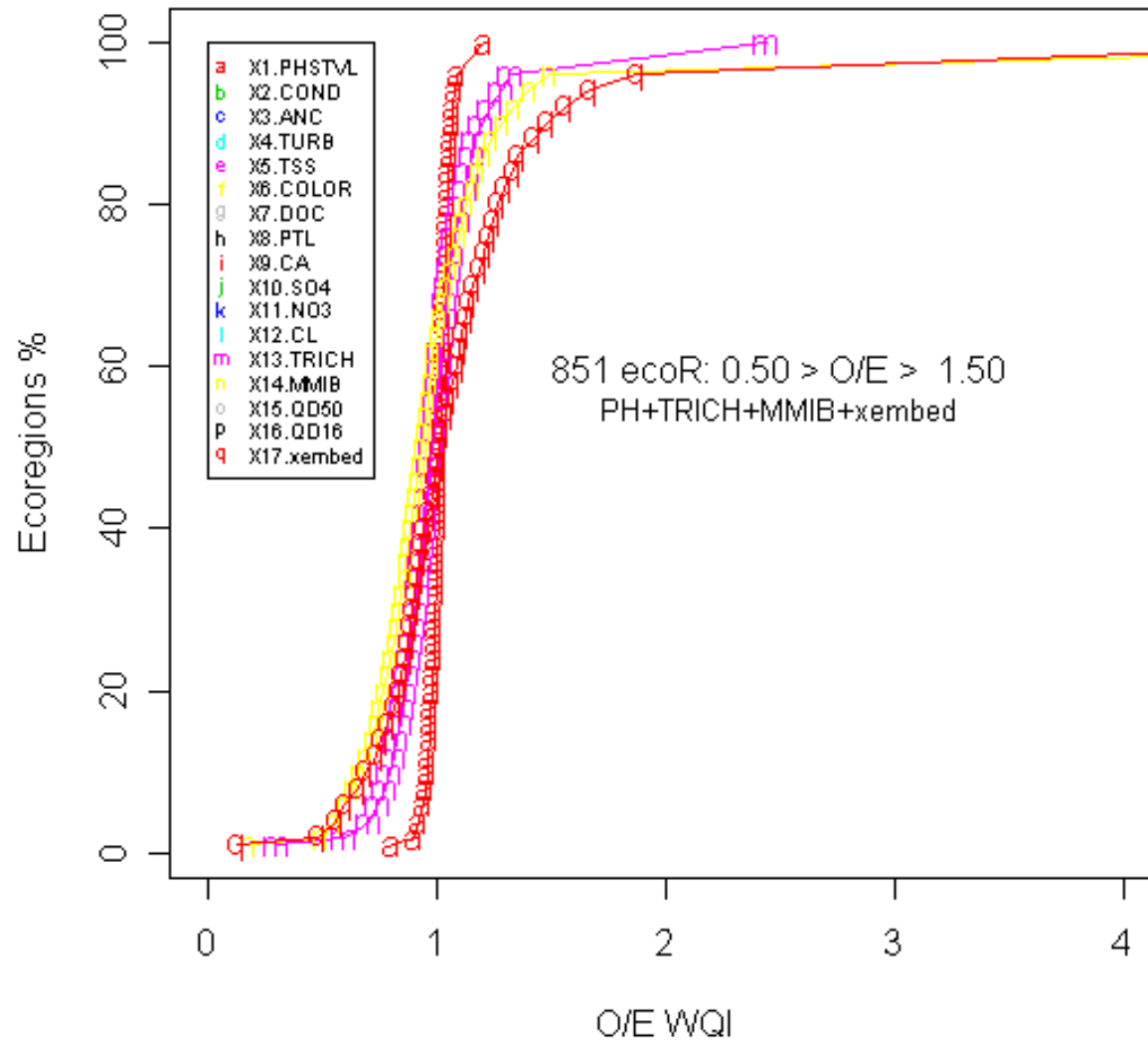
ECOREGIONS (NOV. 2000)

Linking WSA Watersheds and Ecoregions (2)

Expanded quantification using $SL_sqrd(SC+LC)$

- Use a sorted, jointly calculated, list of *squared generalized distance* $SL_sqrd(SC+LC)$ of ecoregions and probability sites, to locate nearest neighbor donors of a mean *WQI* of probability sites (*PRB*) to each ecoregion.
- Use sorted, jointly calculated, list of SL_sqrd of ecoregions and "least-Impacted" reference sites to locate nearest neighbor donors of a mean *WQI* of reference sites (*REF*) to each ecoregion.
- Define each ecoregional *WQI* by the distribution of the ratio of its "*PRB*" and "*REF*" *WQI*, that is, " PRB / REF ". For example:
 - a- "Least Impacted": $0.50 < PRB / REF < 1.50$
 - b- Un-impacted: $PRB / REF < 0.50$
 - c- Impacted: $PRB / REF > 1.50$

Regional Distribution of WQI (grp 2)



Conclusion

- Successful spatial prediction of WSA water quality (*WQI*), using *SC* and *LC* as independent predictors and *WQI* as co-predictors, was demonstrated by 1392 tests, using set-aside watersheds.
- Also successful was replacing set-aside “host” data in 1392 tests with data from the nearest neighbors of each host watershed, which supported transfer of *WQI* data to Ecoregions.

63 Water Quality indicators (W) , Soil Characteristics and land cover (SL) of Watersheds, and 1,392 observations summarized by 63 sample mean $\boldsymbol{\mu}$ and 63x63 covariance $\boldsymbol{\Sigma}$

Watershed 1	$W_{11}, W_{12}, \dots, W_{1,17}, SL_{1,18}, \dots, SL_{1,51} \dots, SL_{1p}$
Watershed 2	$W_{21}, W_{22}, \dots, W_{2,17}, SL_{2,18}, \dots, SL_{2,51} \dots, SL_{2p}$
$p = 17+34+12 = 63, n = 1,392$	
Watershed n	$W_{n1}, W_{n2}, \dots, W_{n,17}, SL_{n,18}, \dots, SL_{n,51}, \dots, SL_{np}$
Mean $\boldsymbol{\mu}$	$\mu_1, \mu_2, \dots, \mu_p$
Variance $\boldsymbol{\sigma}$	$\sigma_{11}, \sigma_{22}, \dots, \sigma_{pp}$
Covariance $\boldsymbol{\Sigma}$	p -dimensional and symmetric, with the diagonal = $\boldsymbol{\sigma}$

Elements of the variance-covariance matrix $\Sigma_{p \times p}$ prior to splitting into *predictor* and *predicted* parts

σ_{11}	σ_{12}			$\sigma_{1,p-1}$	σ_{1p}
σ_{21}	σ_{22}			$\sigma_{2,p-1}$	σ_{2p}
σ_{p1}	σ_{p2}			$\sigma_{p,p-1}$	σ_{pp}

Partitioned $\Sigma_{p \times p}$ with $\Sigma_{1 \times 1}$ along with its covariance Σ_{21} and Σ_{12} isolated from the predictor variables Σ_{22}

$\Sigma_{11} = \sigma_{11}$	$\Sigma_{12} = [\sigma_{12} \cdots \sigma_{1p}]$
Σ_{21}	Predictors: Σ_{22}

Regression with predictor origins= μ , conditional probability, Variance Reduction (VR), and Squared Distance (SL_sqrd)

- $x_1 = \mu_1 + b_2 (x_2 - \mu_2) + \dots + b_p (x_p - \mu_p)$
- x_1 differs from μ_1 by the degree of its correlation with all predictors, thus, $x_1 =$ conditional expectation.
- $E(x_1 \mid x_2 \dots x_p) = \mu_1 + \Sigma_{12} \Sigma_{22}^{-1} (x_{p-1} - \mu_{p-1})$
- Similarly, Conditional Variance is:
- $V_c(x_1 \mid x_2, \dots, x_p) = \sigma_{11} - \Sigma_{12} \Sigma_{22}^{-1} \Sigma_{21}$
- Variance Reduction VR compares V_c with sample variance V_s
- $VR = 100[1.0 - (V_c / V_s)^{1/2}]$
- $SL_sqrd = (\mathbf{x} - \boldsymbol{\mu})' \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu})$