# Fusion of Similarity Measures to Characterize Sample Matrix Effects

**Callan Norby, John H. Kalivas**
Department of Chemistry
Idaho State University
921 S. 8th Ave., STOP 8023 Pocatello, ID 83209, USA
norbcall@isu.edu, kalijohn@isu.edu

**Idaho State UNIVERSITY**

## Abstract

Multivariate calibration applied to spectroscopic data is firmly rooted in the field of analytical chemistry. Over the past several decades, numerous methods have been developed to deduce a calibration model to predict new analyte values with sufficient accuracy and precision. These calibration models produce good results when calibration (primary) and new prediction (secondary) samples are measured under similar conditions. However, inherent sample matrix effects and measurement conditions for the secondary samples are often dissimilar to calibration samples resulting in inaccurate and imprecise predictions. To combat this issue, calibration maintenance by model updating can be used to manipulate the calibration model to adapt to the secondary conditions.

Currently, evaluations of traditional and new calibration maintenance methods by researchers are performed without any consideration for the degree of difference between the primary and secondary data sets. Needed is a method that assesses the degree of difference between primary and secondary data sets for a robust evaluation of any model updating method. In order to solve this problem, multiple similarity measures are utilized in this presentation for a fusion consensus assessment of the degree of difference between the primary and secondary spectra assuming equal distributions of analyte values. Results will be shown for spectral data sets of varying similarity.

## Objective

- Characterize the similarity between two data sets with the same prediction property using 15 similarity measures.
- Validate the method of using similarity measures by correlating the projected similarity to the relative prediction error between data sets.

## Approach

**Part 1: Calculating Indicator of Spectral Uniqueness (ISU)**

Table 1. Similarity measure values for a sample at a single eigenvector window.

| $X_p$ | $X_s$ |
|-------|-------|
| $d_{1p}$ | $d_{1s}$ |
| $d_{2p}$ | $d_{2s}$ |
| $\vdots$ | $\vdots$ |
| $d_{np}$ | $d_{ns}$ |

- For a single sample removed from secondary ($X_s$), all similarity measures are calculated with respect to primary ($X_p$) and the remaining sample spectra in $X_s$
- $d_{ip}$ and $d_{is}$ represent the $i$th primary and secondary similarity measures respectively, where $1 \leq i \leq n$ (integer only) and $n$ is the number of similarity measures

Table 2. Scaled similarity measure values for a sample at a single eigenvector window.

| $X_p$ | $X_s$ |
|-------|-------|
| $d'_{1p}$ | $d'_{1s}$ |
| $d'_{2p}$ | $d'_{2s}$ |
| $\vdots$ | $\vdots$ |
| $d'_{np}$ | $d'_{ns}$ |
| $\dfrac{\sum_{i=1}^{n} s'_{ip}}{n}$ | $\dfrac{\sum_{i=1}^{n} s'_{is}}{n}$ |

- Similarity measures scaled using the equations below

$$d'_{is} = \frac{d_{is}}{d_{ip}+d_{is}} \qquad d'_{ip} = \frac{d_{ip}}{d_{ip}+d_{is}}$$

where $0 \leq d'_{is}, d'_{ip} \leq 1$, and

$$d'_{is} + d'_{ip} = 1$$

- Similarity between the removed sample and the respective space increases as the similarity measure value approaches 0
- Average similarity measure value is then calculated with respect to both spaces
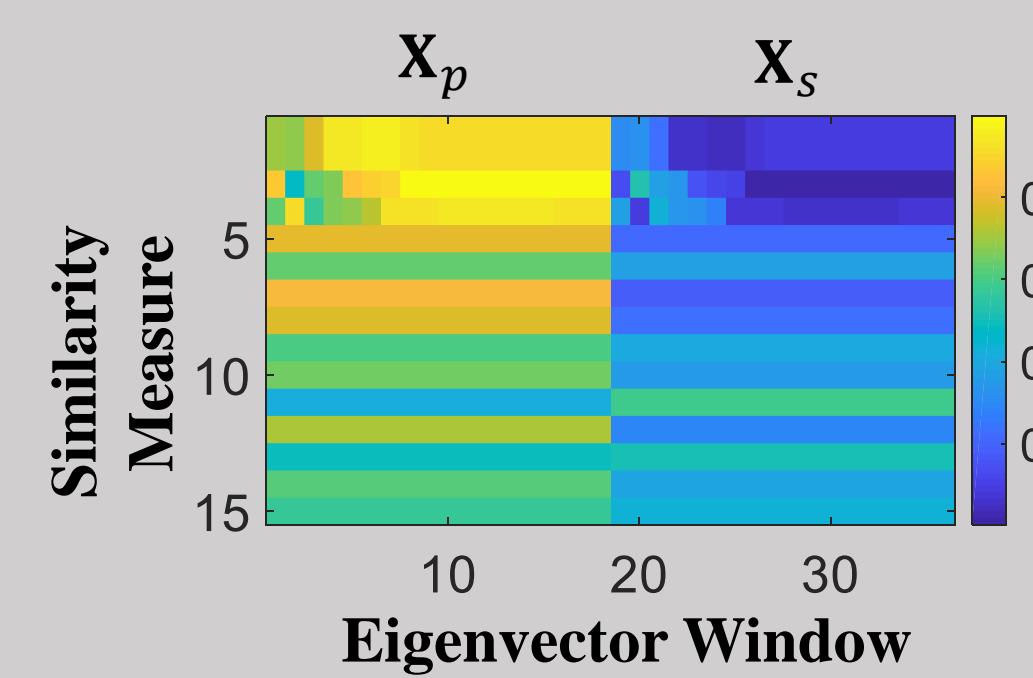
### Similarity Measures (One Sample)

Figure 1. Image of scaled similarity measures for a single sample at 1-rank eigenvector windows.

- Characterizes sample similarity to the primary and secondary data respectively
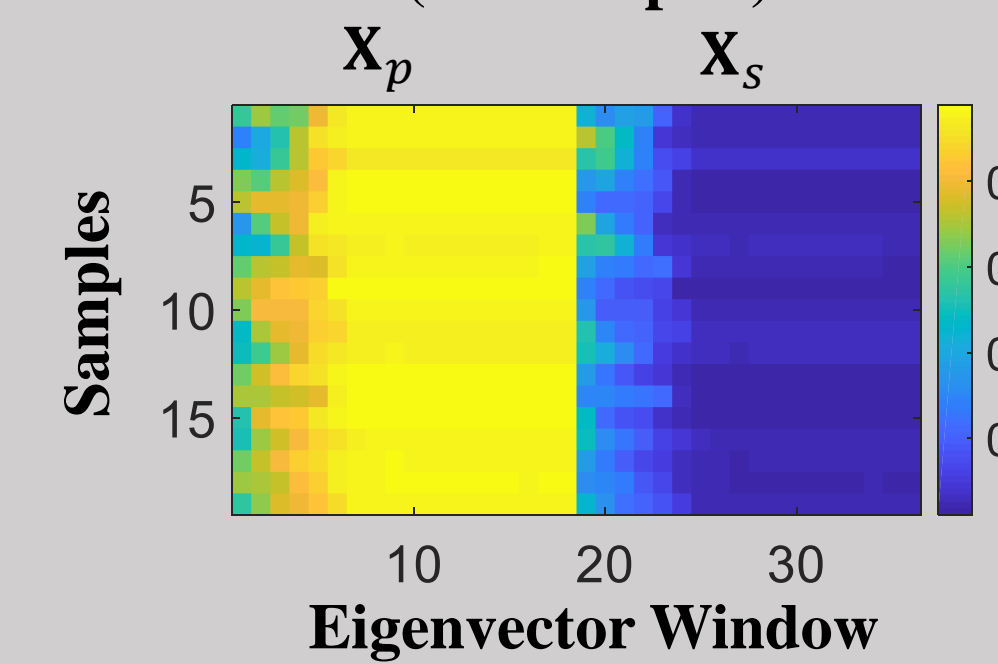
### Similarity Measures (All Samples)

Figure 2. Image of average similarity measure value for every sample in $X_s$ (21) at 1-rank eigenvector windows.

- Procedure repeated for all samples in $X_s$
- Average value across all samples calculated

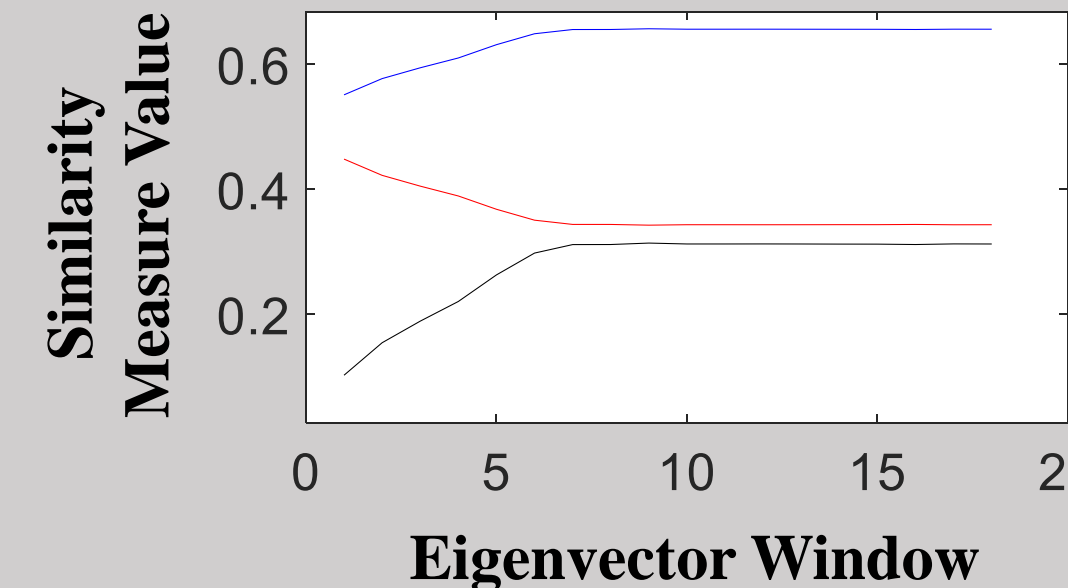### Average Similarity Measure Value (All Samples)

Figure 3. Average similarity measure value across all samples at 1-rank eigenvector windows.

- ISU calculated by subtracting average similarity measure value with respect to secondary from average similarity measure value with respect to primary
- Subtraction performed at last eigenvector window
  - Negates need for model selection

**Part 2: Correlating ISU with Relative Prediction Error**

Build $b_p$   Build $b_{s-i}$

$x_i$

$$\frac{|\hat{y}_{pi} - y_i|}{y_i} \qquad \frac{|\hat{y}_{si} - y_i|}{y_i}$$

- Single sample ($x_i$) removed from $X_s$
- Using PLS1, construct regression coefficients with respect to $X_p$ and $X_{s-i}$
- Predict sample out with respect to $X_p$ and $X_{s-i}$
  - Calculate relative prediction errors
- Repeat for all samples in $X_s$
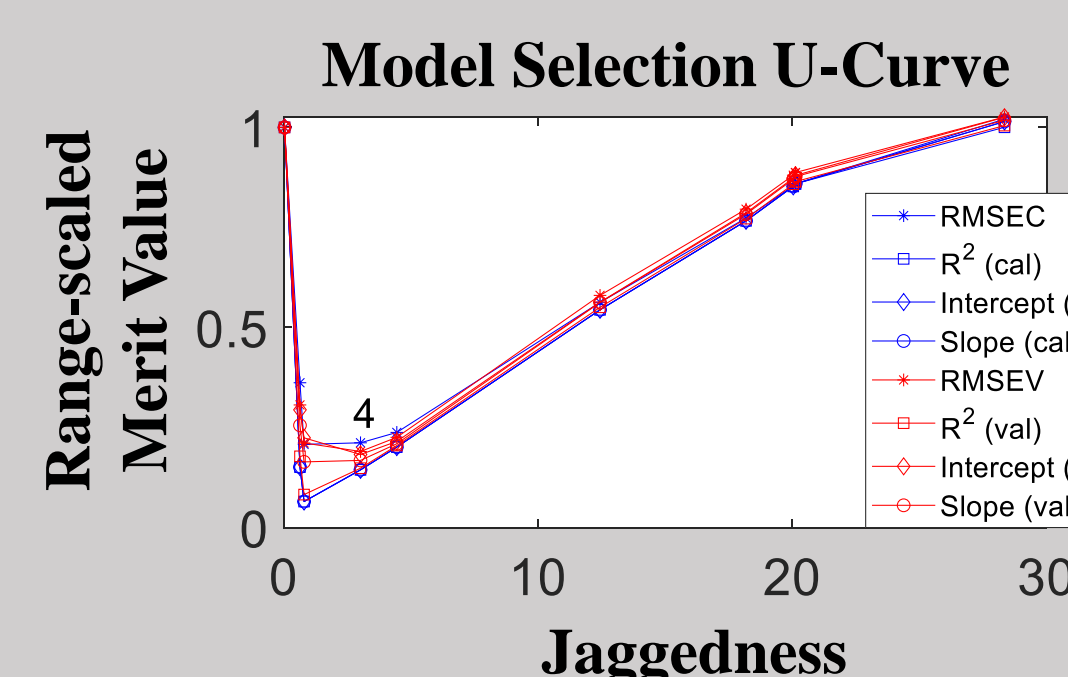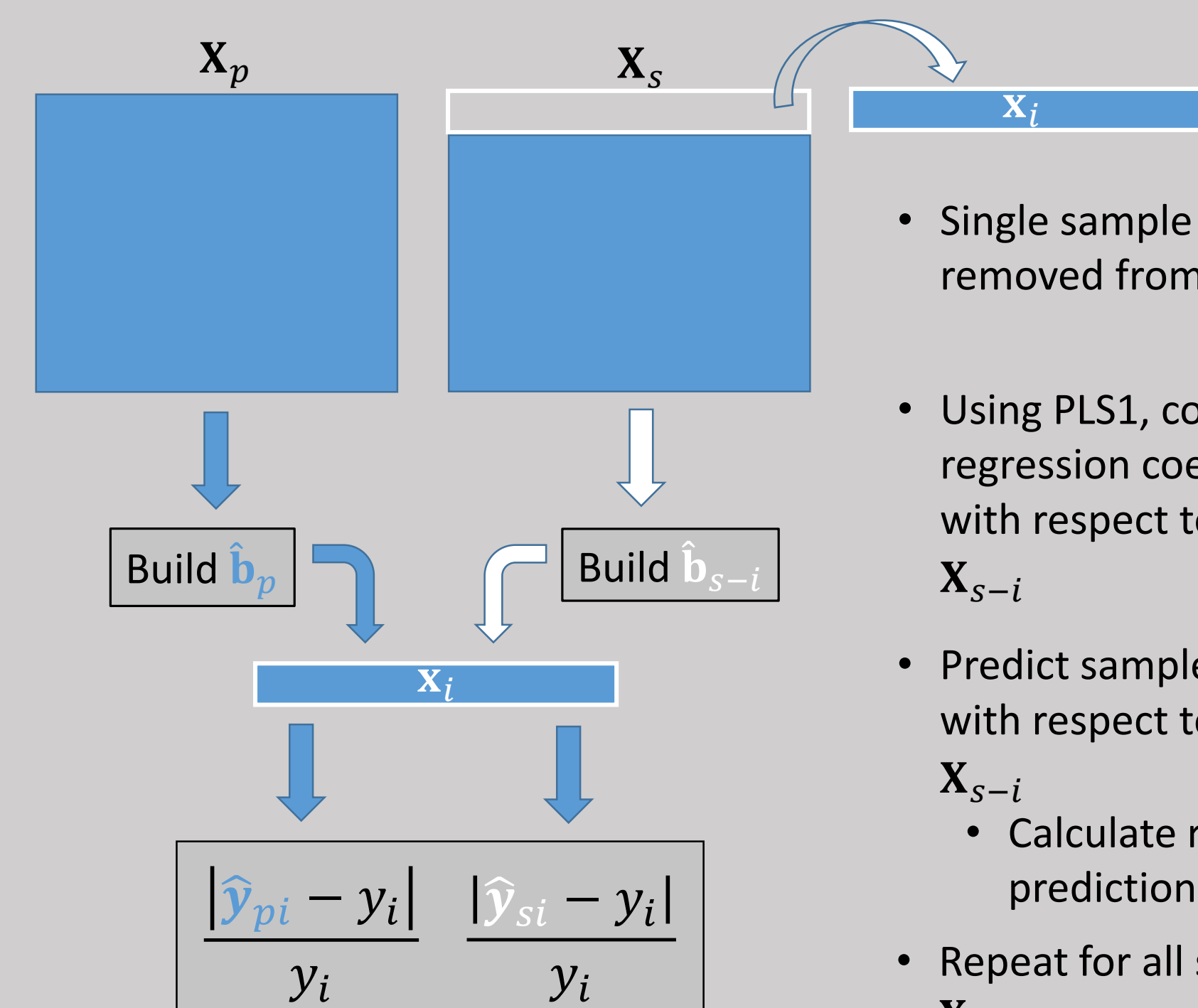
### Model Selection U-Curve

Figure 4. U-curve with all model selection merits at all eigenvector windows

- Select an optimal latent variable (LV) model for $X_p$
- Consensus selection is 4 LV model

### ISU vs. Relative Prediction Error (base corrected, all samples)
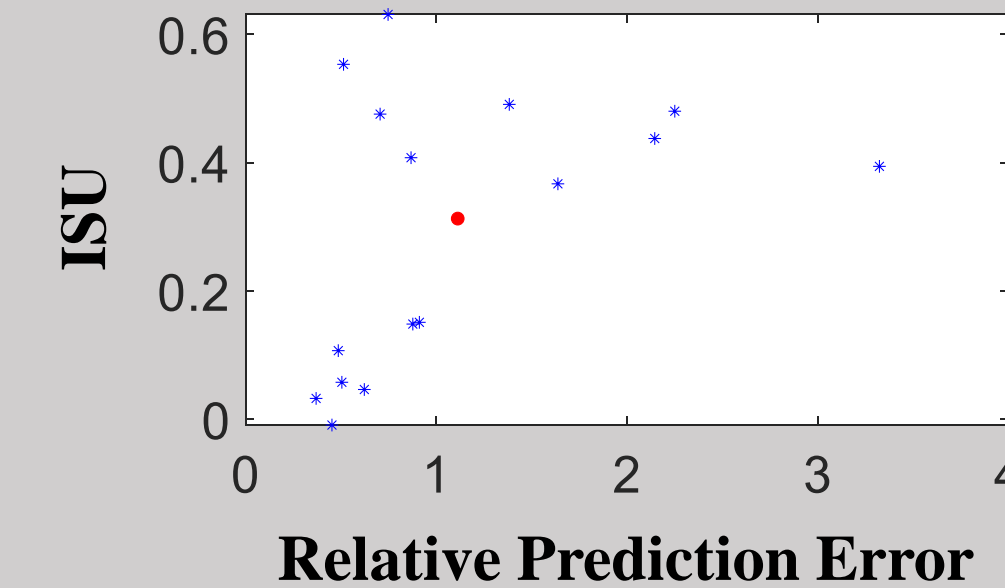
Figure 5. Correlation of relative prediction error to ISU value for all samples / overall average correlation

- Average relative prediction error (base corrected) calculated at selected LV with respect to $X_p$ and $X_{s-i}$

$$\text{Average Relative Prediction Error (base corrected)} = \frac{\sum_{i=1}^{m_p} \frac{|\hat{y}_{pi} - y_i|}{y_i}}{m_p} - \frac{\sum_{i=1}^{m_s} \frac{|\hat{y}_{si} - y_i|}{y_i}}{m_s}$$

- $m_p$ and $m_s$ denote the number of samples in primary and secondary respectively

## Data Set

Temperature
19 samples containing mixtures of ethanol, water, and isopropanol measured at 5 different temperatures via NIR spectrometer.
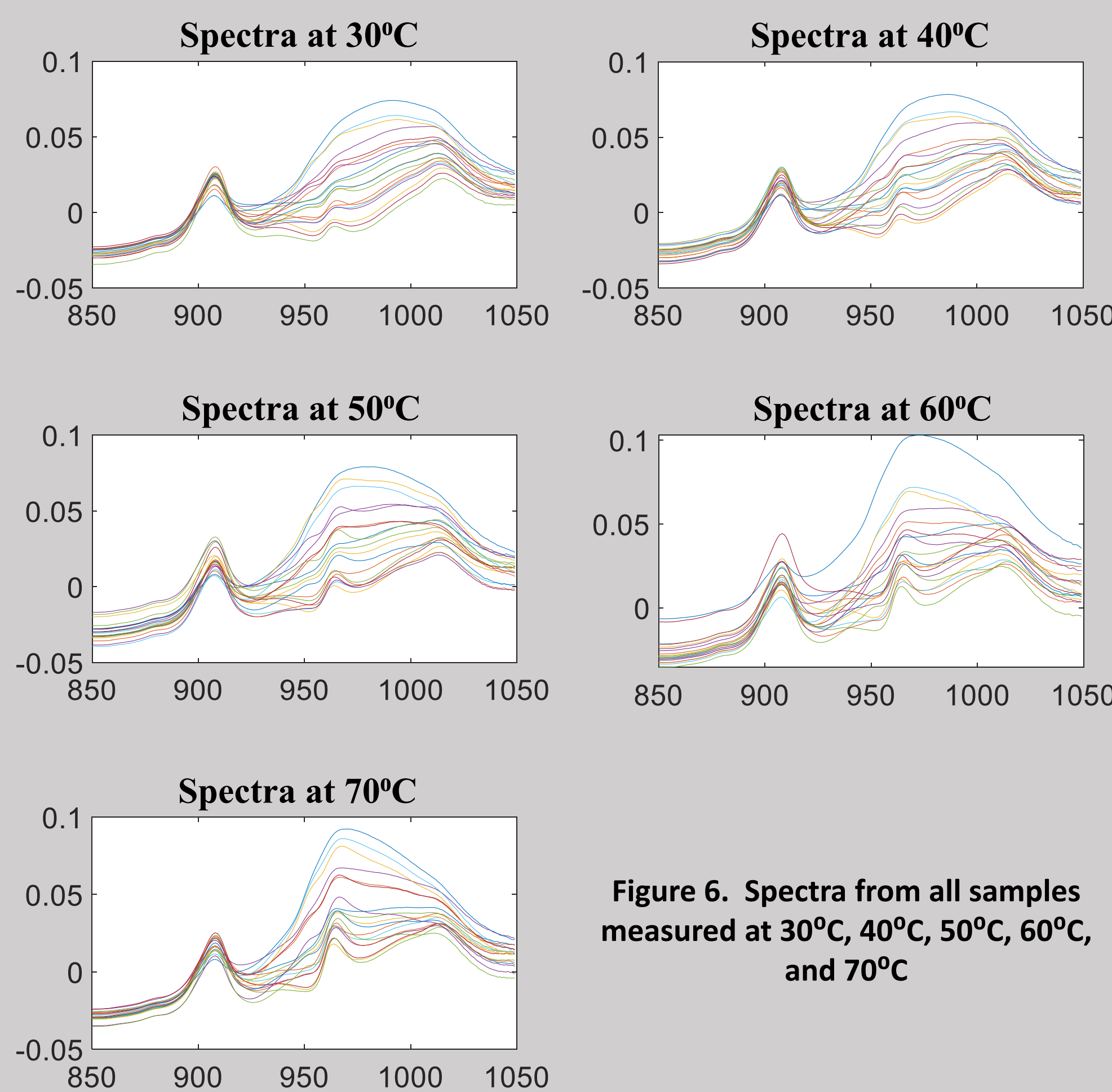
### Spectra at 30°C

### Spectra at 40°C

### Spectra at 50°C

### Spectra at 60°C

### Spectra at 70°C

Figure 6. Spectra from all samples measured at 30°C, 40°C, 50°C, 60°C, and 70°C

## Results

**Temperature Data Set Results**

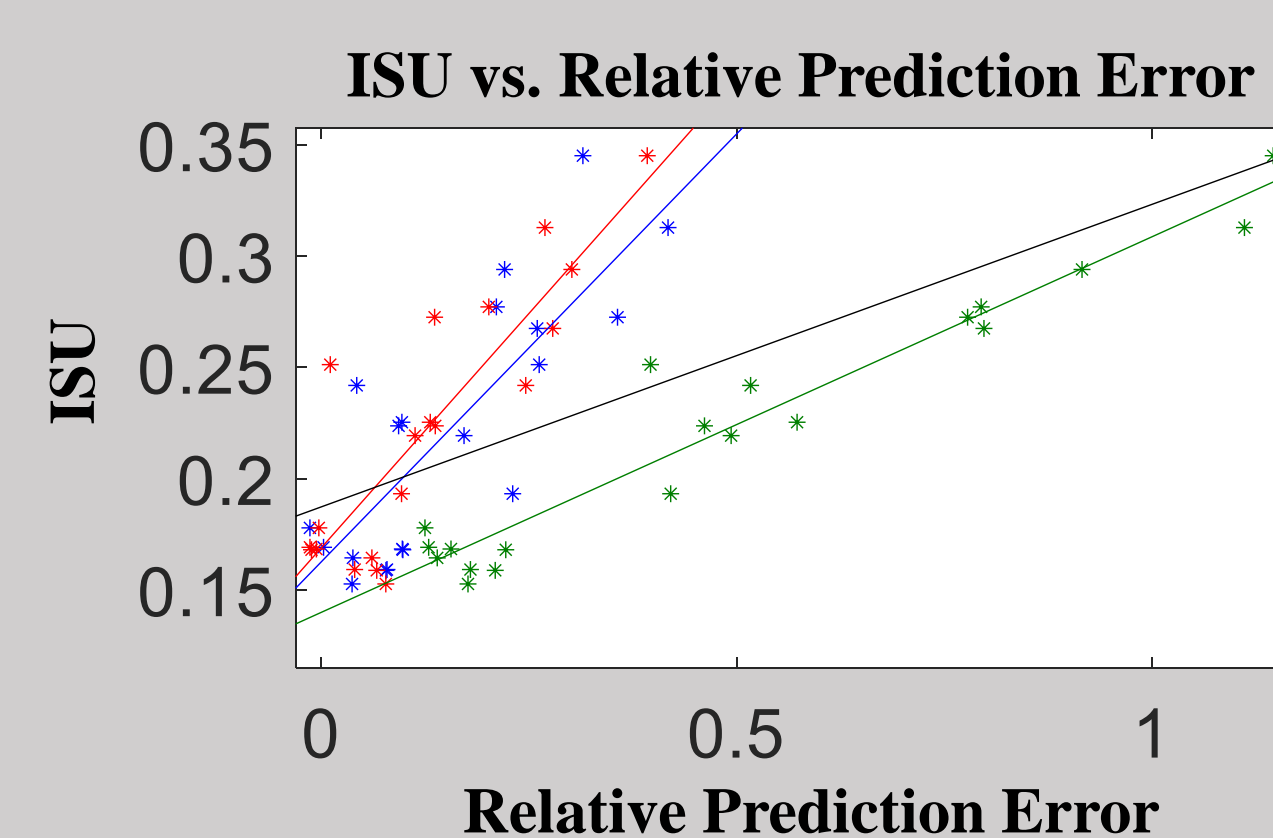### ISU vs. Relative Prediction Error

Figure 7. ISU correlation with relative prediction error for all prediction properties in Temperature

Table 3. ISU vs. relative prediction error correlation values for Temperature data set

| Prediction Property | $R^2$ | Intercept |
|---------------------|-------|-----------|
| Isopropanol | 0.7383 | 0.1688 |
| Water | 0.6565 | 0.1625 |
| Ethanol | 0.9351 | 0.1398 |
| Composite | 0.4078 | 0.1872 |

- Separate trends for each prediction property
- Indicator that analyte information must be accounted for

## Conclusion / Future Work

- ISU criterion is effective at assessing similarity between data sets
  - ISU correlation to prediction error is analyte dependent
    - Account for by including $y$ measures
- Add more $X$ similarity measures
- Evaluate preprocessing methods

## Math Appendix / Similarity Measures

Table 4. Vector-to-space similarity measures with corresponding equations (require a tuning parameter window).

| Mahalanobis Distance | Q-Residual | Sin$\theta$ | Divergence Criterion |
|---------------------|------------|-------------|---------------------|
| $\tilde{C} = \dfrac{\tilde{X}^T \tilde{X}}{n}$ <br> $\tilde{C} = USV^T$ <br> $\tilde{C}_k^+ = U_k S_k^{-1} V_k^T$ <br> $MD_i = \sqrt{(x_i - \bar{x})^T \tilde{C}_k^+ (x_i - \bar{x})}$ | $x_i^\perp = (I - V_k V_k^T)x_i$ <br> $X = USV^T$ <br> $Q_i = \|x_i^\perp\|$ | $\sin\theta_i = \dfrac{\|x_i^\perp\|}{\|x_i\|}$ | $DC_i = \left\| \frac{1}{2} tr((X_i - C)(X_i^+ - C_i^+)) + \frac{1}{2} tr((X_i^+ - C_i^+)(x_i - \bar{x})(x_i - \bar{x})^T) \right\|$ |

Table 5a. Vector-to-vector similarity measures with corresponding equations.

| Procrustes Analysis (unconstrained) | Procrustes Analysis (Constrained) | Extended Inverted Signal Correction Difference |
|-------------------------------------|-----------------------------------|-----------------------------------------------|
| $\bar{X} = \rho \tilde{X} H$ <br> $\bar{X}^T \bar{X} = USV^T$ <br> $H = UV^T$ <br> $\rho = \dfrac{tr(S)}{tr(\tilde{X}\tilde{X}^T)}$ | $\tilde{X} = \rho_i X_i H_i$ <br> $X_i^T \tilde{X} = U_i S_i V_i^T$ <br> $H_i = U_i V_i^T$ <br> $\rho_i = \dfrac{tr(S_i)}{tr(X_i \tilde{X}^T)}$ | $\bar{X} = X_i T_i$ <br> $\tilde{X} = \bar{X} T$ <br> $T_i = X_i^+ \tilde{X}$ <br> $T = \bar{X}^+ \tilde{X}$ | $\bar{x} = x_i + X_c b$ <br> $d = \bar{x} - x_i = X_c b$ <br> $X_c = \left(1, x_i, x_i^2, \dfrac{d}{d\lambda}x_i, \dfrac{d^2}{d\lambda^2}x_i, \lambda, \lambda^2, \ln(\lambda)\right)$ |
| $\Delta H_i = \|H_i - H\|_F$ <br> $\Delta\rho_i = |\rho_i - \rho|$ | $\Delta T_i = \|T_i - T\|_F$ | $EISCD_b = \|b\|$ <br> $EISCD_{X,b} = \|X_c b\|$ |

Table 5b. Vector-to-vector similarity measures with corresponding equations.

| Determinant | Inner Product Correlation | Euclidian Distance | $1 - \cos\theta$ |
|-------------|---------------------------|--------------------|------------------|
| $det_i = Det\left(\begin{bmatrix} x_i^T \\ \bar{x}^T \end{bmatrix}\begin{bmatrix} x_i & \bar{x} \end{bmatrix}\right)$ <br> $= (\|x_i\|\|\bar{x}\|\sin\theta_i)^2$ | $1 - r_i = 1 - \dfrac{tr(X_i^T \bar{X})}{\sqrt{tr(X_i^T X_i)tr(\bar{X}^T \bar{X})}}$ | $d_i = \|x_i - \bar{x}\|$ | $1 - \cos\theta_i = 1 - \dfrac{x_i^T \bar{x}}{\|x_i\|\|\bar{x}\|}$ |

**Notation for Table 1, 2a, 2b**

- $k$ subscript denotes the arbitrary number of eigenvectors or latent variables selected
- $\bar{x}$ is the column-wise mean vector of $X$
- Outer product arrays $\bar{X}$ and $X_i$ are computed by $\bar{X} = \bar{x}\bar{x}^T$ and $X_i = x_i x_i^T$
- $\| \ \|_F$ denotes the Frobenius norm
- $\lambda$ is the vector of wavelengths
- Four EISCD similarity measures are created by swapping $x_i$ and $\bar{x}$