# ILS to CLS: Synergistic Regression Modeling for Improved Control and Interpretability

## Neal B Gallagher

## Eigenvecter Research, Inc.

Booth 72
Gallagher, SciX, Sep 18-23, 2016

EIGENVECTOR RESEARCH INCORPORATED

---

ILS to CLS: Synergistic Regression Modeling for Improved Control and Interpretability
Neal B Gallagher, Eigenvecter Research, Inc.

Inverse least squares (ILS) models such as partial least squares and principle components regression are popular regression tools for chemometrics modeling. A major reason for this popularity is that extensive infrastructure has been developed to make model identification fast and easy. Additionally, statistical diagnostics provide tools to develop useful models for exploratory analysis and quantification tasks. However, although much work has gone into developing tools for interpretation of ILS models, classical least squares (CLS) models are superior for interpretation. CLS also provides more control during model identification *and* application because the model form is amenable to constraints that incorporate known physics and chemistry. Unfortunately, identification of CLS models is often more difficult than for ILS models – a property often attributed to interference signal present in measurements. This talk will show that the advantages of ILS and CLS can be used synergistically resulting in models that provide enhanced diagnostics and interpretability. Two examples typically modeled using ILS will be shown.

Booth 72
Gallagher, SciX, Sep 18-23, 2016

EIGENVECTOR RESEARCH INCORPORATED

# Summary

- Demonstrate the theory and practice of using inverse least squares (PLS and PCR) with classical least squares methods (CLS and ELS)
- PLS & PCR are fast and easy to identify
  - difficult to control, difficult to interpret and can be confused by coincidental correlation
- CLS & ELS can be difficult to identify
  - Extended Least Squares (extended mixture model)
  - easy to add what we know via constraints, easy to interpret shows if coincidental correlation is present

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
RESEARCH INCORPORATED

---

# Inverse Least Squares

$$\mathbf{y} = \mathbf{Xb}$$

Partial Least Squares (PLS) and
Principal Components Regression (PCR)

$\mathbf{X}$ is the predictor (e.g., measured spectra)

$\mathbf{y}$ is the predictand (e.g., concentration, univariate)

$\mathbf{b}$ is the regression vector (univariate)

$$\mathbf{Y} = \mathbf{XB}$$ Multivariate $\mathbf{Y}$ is the more general case.

$$\mathbf{B} = \mathbf{W}\left(\mathbf{T}^T\mathbf{T}\right)^{-1}\mathbf{T}^T\mathbf{Y}$$

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
RESEARCH INCORPORATED

# ILS to ELS

$$\mathbf{Y} = \mathbf{XB} \qquad \mathbf{T} = \mathbf{XW}$$

$\mathbf{T}$ are scores ($\mathbf{Y}$ is in column space of $\mathbf{T}$).
$\mathbf{W}$ are weights for PLS and **loadings** for PCR
($\mathbf{B}$ is in the column space of $\mathbf{W}$).

… rearrange the variance...

$$\begin{bmatrix} \mathbf{Y} & \mathbf{T}_\perp \end{bmatrix} = \mathbf{X} \begin{bmatrix} \mathbf{B} & \mathbf{W}_\perp \end{bmatrix} \qquad \text{OPLS}$$

$$\mathbf{X} = \begin{bmatrix} \mathbf{Y} & \mathbf{T} \end{bmatrix} \begin{bmatrix} \mathbf{S} & \mathbf{P} \end{bmatrix}^T \qquad \text{ELS}$$

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR** RESEARCH INCORPORATED

# ELS Objective Function

$$O\left(\mathbf{Y}, \mathbf{T}_\perp, \mathbf{S}, \mathbf{P}_\perp\right) =$$

$$\left(\mathbf{X} - \begin{bmatrix} \mathbf{Y} & \mathbf{T}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{S} & \mathbf{P}_\perp \end{bmatrix}^T \right)^T \left(\mathbf{X} - \begin{bmatrix} \mathbf{Y} & \mathbf{T}_\perp \end{bmatrix} \begin{bmatrix} \mathbf{S} & \mathbf{P}_\perp \end{bmatrix}^T \right)$$

$$+ \left(\mathbf{I} - \mathbf{P}_\perp\right)^T \mathbf{A} \left(\mathbf{I} - \mathbf{P}_\perp\right)$$

$\mathbf{A}$ is a diagonal penalty factor. Orthogonality condition
on $\mathbf{P}$ retains good mathematical conditioning.

Very good estimates of $\mathbf{Y}$ and $\mathbf{T}$ are available from PCR.

$\mathbf{Y}$ and $\mathbf{S}$ are often non-negative.

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR** RESEARCH INCORPORATED

# Why is this Important?

- Inverse least squares methods like PCR and PLS are fast and easy to identify
  - Infrastructure and statistics are well defined
  - Interpretability can be difficult (**B** *are not* spectra!)
  - Many constraints for **B** don't make physical sense
    - non-negativity and smoothness aren't generally applicable
    - this hampers including physical knowledge into the objective function
  - ILS provides very good guesses for **Y** and **T**
    - initial solutions for ELS

$$\begin{bmatrix} \mathbf{Y} & \mathbf{T}_{\perp} \end{bmatrix} = \mathbf{X} \begin{bmatrix} \mathbf{B} & \mathbf{W}_{\perp} \end{bmatrix}$$

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
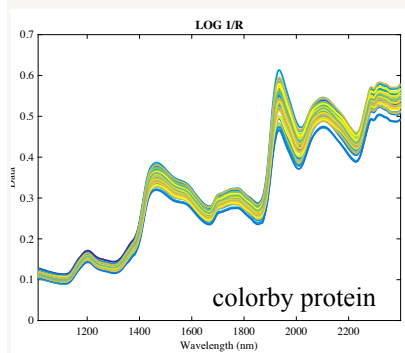RESEARCH INCORPORATED

# Why is this Important

- Forward least squares methods like CLS and ELS can be more difficult to identify
  - Infrastructure is less well developed
  - Can include more of what we know via constraints during model identification *and* application
  - Interpretability is as good as it gets (**S** *are* spectra!)
    - allows iterative model identification because the identification process teaches about the problem

$$\mathbf{X} = \begin{bmatrix} \mathbf{Y} & \mathbf{T} \end{bmatrix} \begin{bmatrix} \mathbf{S} & \mathbf{P} \end{bmatrix}^{T}$$

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
RESEARCH INCORPORATED

# NIR: Wheat Protein and Moisture
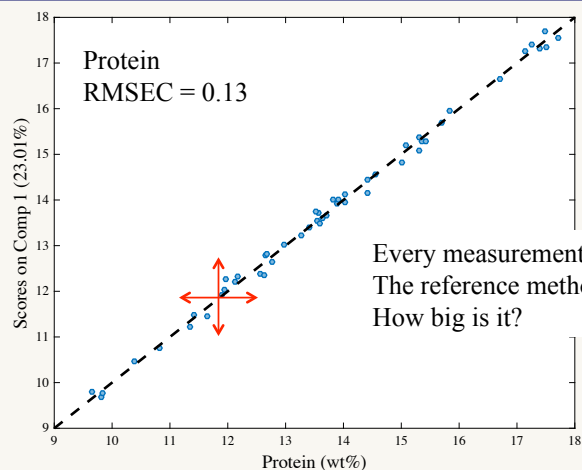
LOG 1/R

colorby protein

Info courtesy D. Hopkins.
Data courtesy P. Williams and K. Norris

- Hard red winter wheat ground
- calibration and validation sets measured at different times
- Cary-14 spectrometer system, 1000 – 2598.4 nm at 1.6 nm intervals, 3 nm resolution
- Protein by Kjeldahl, each sample measured 16 times, averaged
  - Estimated Standard Error of Laboratory, 0.14% protein for the averaged results
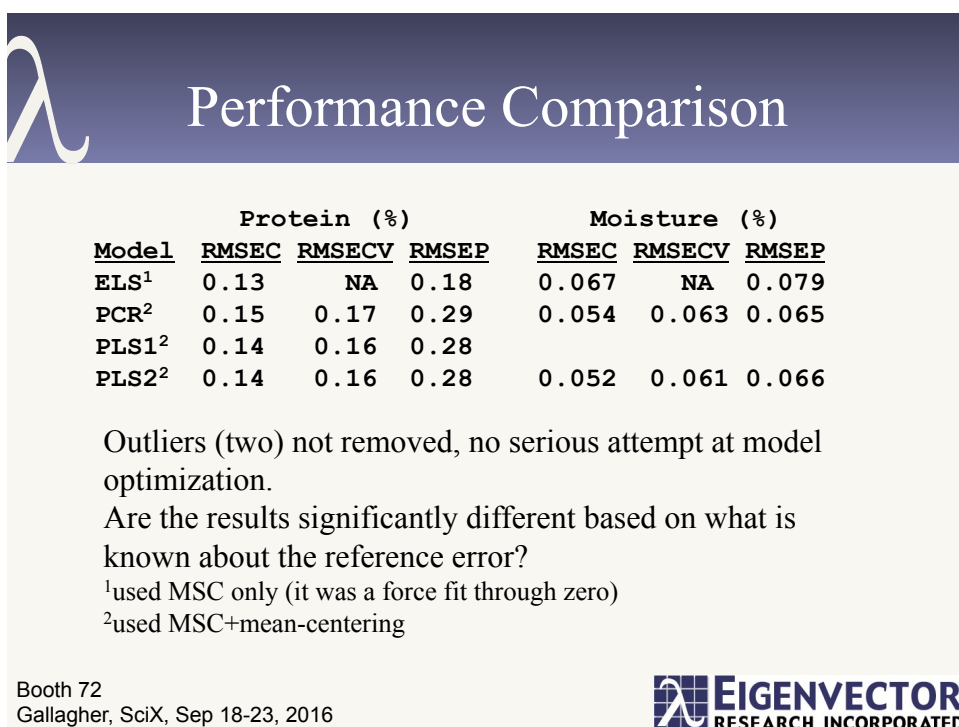
Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
RESEARCH INCORPORATED

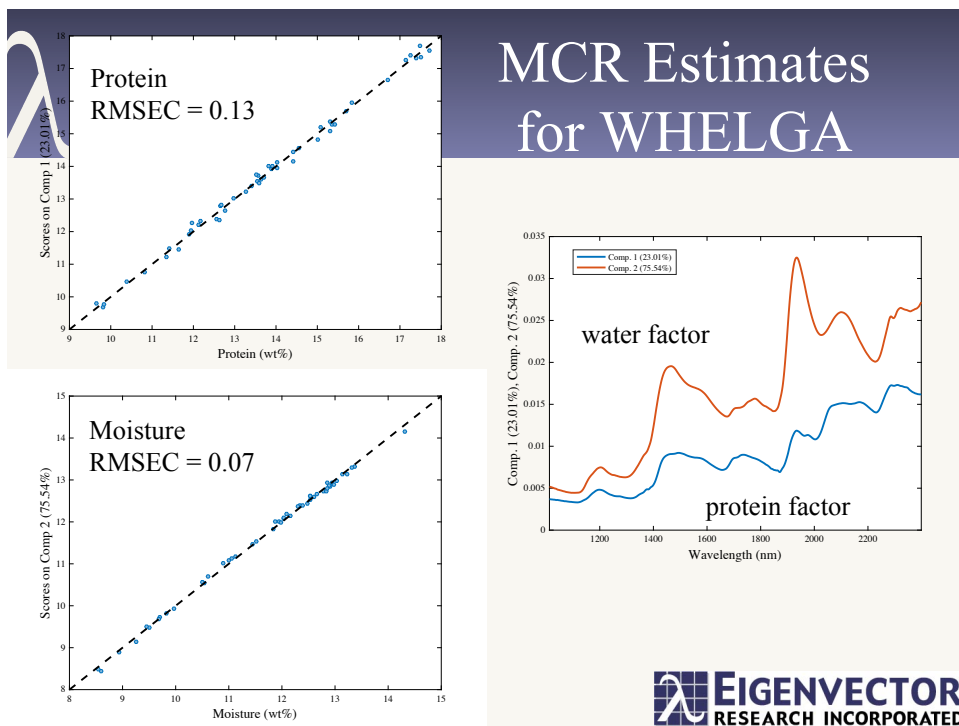# Error in the Reference Method

Protein
RMSEC = 0.13

Every measurement has an error.
The reference method has an error.
How big is it?

Protein (wt%)

Scores on Comp 1 (23.01%)

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
RESEARCH INCORPORATED

# Scores and Weights



The regression vector, **b**, is a linear combination of the weights, $\mathbf{W}_K$.

Booth 72
Gallagher, SciX, Sep 18-23, 2016

# Performance Comparison

| Model | Protein (%) | | | Moisture (%) | | |
|---|---|---|---|---|---|---|
| | RMSEC | RMSECV | RMSEP | RMSEC | RMSECV | RMSEP |
| ELS[1] | 0.13 | NA | 0.18 | 0.067 | NA | 0.079 |
| PCR[2] | 0.15 | 0.17 | 0.29 | 0.054 | 0.063 | 0.065 |
| PLS1[2] | 0.14 | 0.16 | 0.28 | | | |
| PLS2[2] | 0.14 | 0.16 | 0.28 | 0.052 | 0.061 | 0.066 |

Outliers (two) not removed, no serious attempt at model optimization.
Are the results significantly different based on what is known about the reference error?
[1]used MSC only (it was a force fit through zero)
[2]used MSC+mean-centering

Booth 72
Gallagher, SciX, Sep 18-23, 2016

## MCR Estimates for WHELGA



Protein
RMSEC = 0.13

Moisture
RMSEC = 0.07

water factor

protein factor

EIGENVECTOR
RESEARCH INCORPORATED

## Interpretation



water factor

protein factor

first derivative

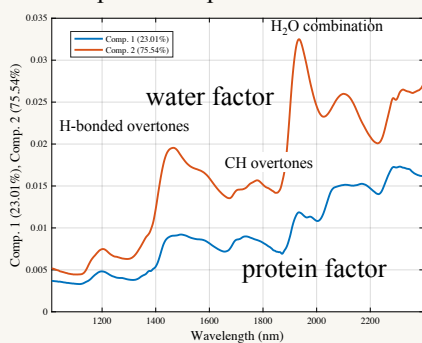PLS regression vector for protein
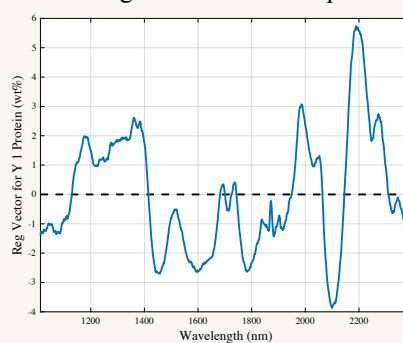
EIGENVECTOR
RESEARCH INCORPORATED

## Interpretation

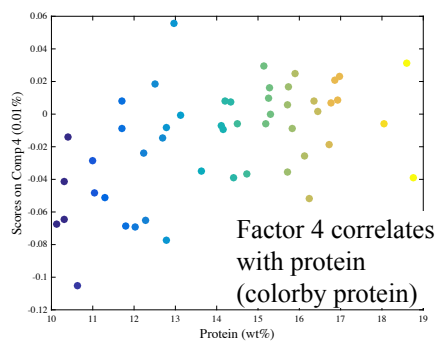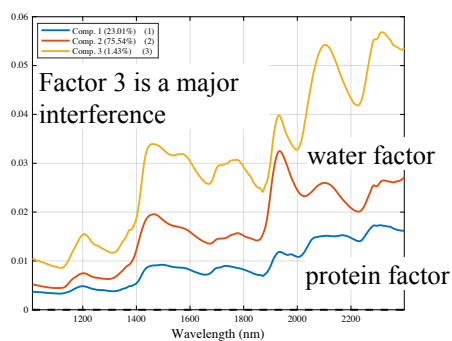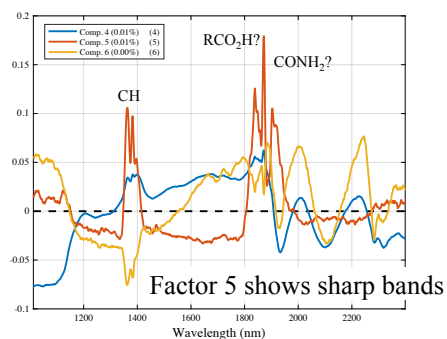ELS spectra for protein and moisture

PLS regression vector for protein

ELS spectra identified using multivariate
curve resolution with soft equality constraints

Booth 72
Gallagher, SciX, Sep 18-23, 2016

EIGENVECTOR
RESEARCH INCORPORATED



## Interferences

Interferences can be interpreted too…
and the estimate of the factors might be
tuned up with additional constraints

Factor 5 shows sharp bands

Factor 3 is a major interference

water factor

protein factor

Factor 4 correlates
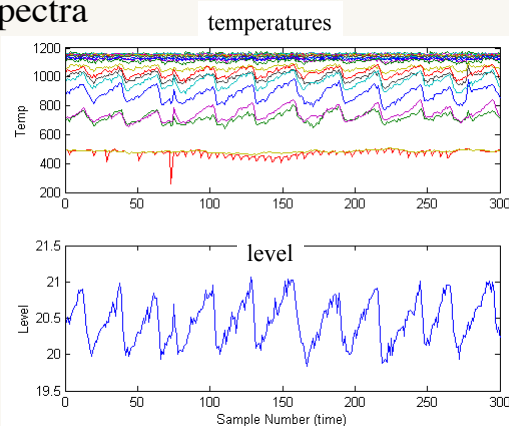with protein
(colorby protein)

# SFCM Data Example

▨ Estimate level in a slurry fed ceramic melter[*]

- ▪ measurements are not spectra

- • measured 20 temperatures (thermocouples) in two vertical thermal wells

- • thermocouples near the surface vary with level

[*]Wise BM, Gallgher NB, Bro R, Shaver JM, Windig W, Koch RS, PLS_Toolbox 3.5, for Use with MATLAB™, Eigenvector Research, Inc. 2004.
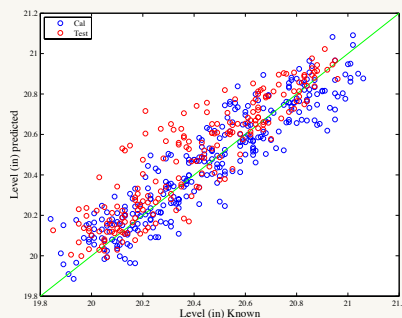
temperatures



Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR RESEARCH INCORPORATED**

---

# Performance Comparison

ELS Calibration and Prediction



Demonstration of ELS/CLS for engineering variables.
Mean-Centering for both PLS and ELS
3 factors for PLS and ELS

|  | Level (in) | |
|---|---|---|
|  | PLS | ELS |
| RMSEC | 0.106 | 0.114 |
| RMSECV | 0.113 | 0.118 |
| RMSEP | 0.138 | 0.145 |

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR RESEARCH INCORPORATED**

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
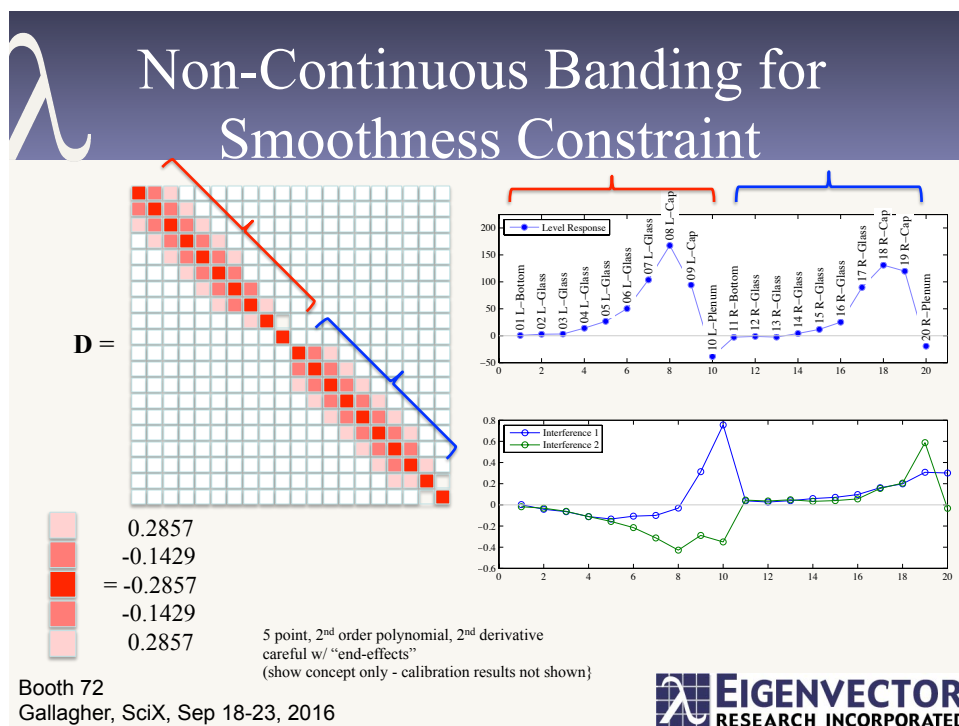RESEARCH INCORPORATED

# Constraints

- Constraints can be employed on both **C** and **S** during ELS model identification
  - e.g., non-negativity, smoothness, priors, time-series lagging, etc…
- … and on **C** during model application
- Allows imposing chemical and physical knowledge into the model identification

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR**
RESEARCH INCORPORATED

# Non-Continuous Banding for Smoothness Constraint

$D =$

| | |
|---|---|
| | 0.2857 |
| | -0.1429 |
| = | -0.2857 |
| | -0.1429 |
| | 0.2857 |

5 point, 2nd order polynomial, 2nd derivative
careful w/ "end-effects"
(show concept only - calibration results not shown}

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR** RESEARCH INCORPORATED

---

# Conclusions

ILS or CLS ?          ILS + CLS !

- Inverse least squares methods like PCR and PLS are fast and easy to identify
  - Interpretability can be difficult (**B** *are not* spectra!)
- Forward least squares methods like CLS and ELS allow more control over model identification
  - Interpretability is as good as it gets (**S** *are* spectra!)

Booth 72
Gallagher, SciX, Sep 18-23, 2016

**EIGENVECTOR** RESEARCH INCORPORATED