# Putting Chemistry Back Into NIR Calibration Models with Gray Classical Least Squares

**Barry M. Wise[1], Donal O'Sullivan[1] and Rasmus Bro[2]**
**[1]Eigenvector Research, Inc.**
**[2]University of Copenhagen**

**EIGENVECTOR RESEARCH INCORPORATED**

# Abstract

There is a resurgence in the use of Classical Least Squares (CLS) models primarily due to their interpretability. When used with spectroscopic systems that follow the Lambert-Beer law CLS models follow naturally from first principles. Unfortunately, CLS models typically do not have the predictive ability of inverse least squares (ILS) models such as Partial Least Squares (PLS) regression: the prediction error of CLS models is usually higher, and often notably so. This is largely due to non-idealities in the data of interest along with the presence of unaccounted for minor components, *e.g.* scatter and baseline variations. PLS models handle these situations by adding components to the model that keep the resulting regression vector orthogonal to the non-ideal variations. In this work we examine a method for developing CLS models with predictive properties competitive with ILS formulations. This is done by using the CLS model "half-residuals" to develop pre-filters with Generalized Least Squares Weighting (GLSW) or External Parameter Orthogonalization (EPO). The result is calibration models that have chemically meaningful estimates of the pure component spectra, interpretable factors for non-idealities and minor components and good predictive ability. Gray CLS models are demonstrated with several NIR data sets and their performance is shown comparable to PLS models.

EIGENVECTOR
RESEARCH INCORPORATED

2

# Outline

- The Classical Least Squares Model
  - Ways to generate spectral residuals
- Clutter Filters
  - EPO: External Parameter Orthogonalization
  - GLSW: Generalized Least Squares Weighting
- Example Data Sets
- Results
  - Choosing the Meta-parameters
  - Diagnostics
- Conclusions

# Classical Least Squares

- CLS often used to develop spectroscopic calibration models
- The CLS assumes the data can be modeled as

$$X = CS^T + E$$

where:

- $X$ $(M \times N_x)$    is the measured spectral response
- $S$ $(N_x \times K)$    is a matrix of pure spectral responses,
- $C$ $(M \times K)$    is a matrix of concentrations and
- $E$ $(M \times N_x)$    is noise or an error matrix.

4

**EIGENVECTOR**
RESEARCH INCORPORATED

# Using the CLS Model

- If $S$ is known, the $\widehat{C}$ can be estimated from
  - $\widehat{C} = XS(S^T S)^{-1}$
- If $S$ is not known, it can be estimated from a (properly designed) calibration data set
  - $\hat{S} = (C^T C)^{-1} C^T X$
- This is often the best way to estimate $S$
  - Models $S$ in the relevant sample matrix
  - Temperature, pressure, scattering effects, etc.

EIGENVECTOR
RESEARCH INCORPORATED

# CLS Spectral Residuals

- Given $\hat{S}$ there are two ways to get spectral residuals
- Conventional $\boldsymbol{R_c}$, estimate $\widehat{\boldsymbol{C}}$ as above then
  - $\boldsymbol{R_c} = \boldsymbol{X} - \widehat{\boldsymbol{C}}\hat{\boldsymbol{S}}^T$
- Half residuals $\boldsymbol{R_h}$, use original $\boldsymbol{C}$ instead of $\widehat{\boldsymbol{C}}$ $\longleftarrow$ $\mathcal{CRUX}$
  - $\boldsymbol{R_h} = \boldsymbol{X} - \boldsymbol{C}\hat{\boldsymbol{S}}^T$
- Note that $\boldsymbol{R_c}$ is orthogonal to $\hat{\boldsymbol{S}}$, whereas $\boldsymbol{R_h}$ is not
  - <u>This will be important later!</u>

EIGENVECTOR
RESEARCH INCORPORATED

# Main Problem with CLS

- As originally formulated, typically not competitive with PLS and other Inverse Least Squares (ILS) approaches on prediction error

- Can't use with un-quantified unknown components

- Factor based methods (PLS, PCR) compensate for non-idealities by going beyond the number of known components

- With CLS you're stuck. ??

EIGENVECTOR
RESEARCH INCORPORATED

# Clutter Orthogonalization Filters

- Typically used as preprocessing in PLS or other ILS models
- Mitigate effect of large variations in spectra *not* related to property of interest
- Consider two popular orthogonalization filters here
  - External Parameter Orthogonalization (EPO)
  - Generalized Least Squares Weighting (GLSW)

EIGENVECTOR
RESEARCH INCORPORATED

# EPO Filter

- Given a matrix $Z$ which represents extraneous variation (matrix effects, clutter), decompose $Z$ as

  - $Z = USV^T$ <span style="color:red">$Z = ?$</span>

- The number of filter factors $k$ must be specified, then

  - $F_{epo} = I - V_k V_k^T$

- $F_{epo}$ is applied to $X$ before calibration (and during prediction), and removes variations in the first $k$ dimensions represented in $Z$

- Equivalent to Extended Mixture Model (EMM)

**EIGENVECTOR** RESEARCH INCORPORATED

# GLSW Filter

- Similar to EPO filter except it shrinks dimensions rather than completely eliminating them

- Starting from the decomposition of $\boldsymbol{Z}$ above then

  - $\boldsymbol{F_{glsw}} = \boldsymbol{VD^{-1}V^T}$

- Where the diagonal elements of $\boldsymbol{D}$ are calculated as

  - $d_i = \sqrt{\left(\frac{s_i^2}{g^2}\right) + 1}$

- Where $s_i$ is the i[th] diagonal element of $\boldsymbol{S}$ and $g$ is a tune-able parameter which controls the shrinkage

EIGENVECTOR
RESEARCH INCORPORATED

# Comparison of EPO & GLSW

**Singular Values of Clutter Source**

**Corresponding Value in GLSW or EPO Filter**

EPO 5 factors

GLS decreasing $g$

EIGENVECTOR
RESEARCH INCORPORATED

# Combining CLS & Filters

- Residuals from CLS models can be used as an estimate of the clutter, $Z$. However,
  - Filter based on $R_c$ has *no effect* as it is orthogonal to $\hat{S}$.
  - $R_h$, on the other hand, contains information about clutter *not* orthogonal to $\hat{S}$.
- Filter based on $R_h$ mitigates clutter *not* orthogonal to $\hat{S}$ that would otherwise lead to additional error in $\hat{C}$.

*CRUX*

**EIGENVECTOR**
RESEARCH INCORPORATED

12

# The CLS Gray Model

- We refer to this combination of a CLS model with a filter based on the half residuals $R_h$ as a "gray model"

- Incorporates aspects of both CLS and ILS models.
  - Based on a first principles model, the CLS "white" part
  - Includes tunable empirical part, EPO or GLSW filter "black" part
  - This model has a single adjustable parameter ($k$ or $g$)

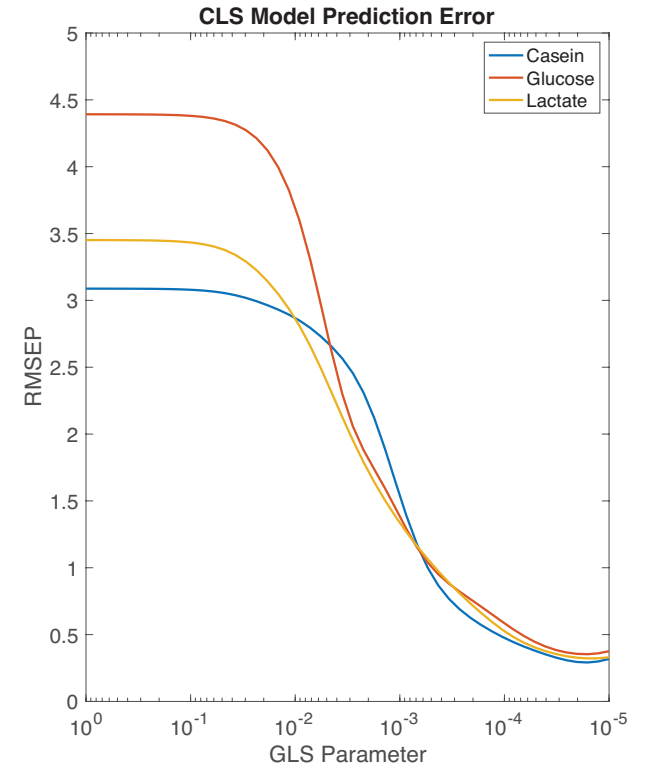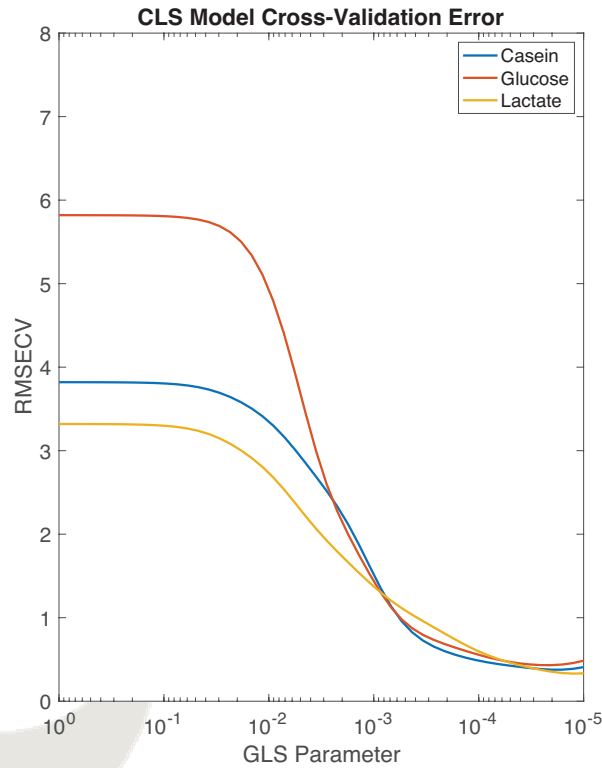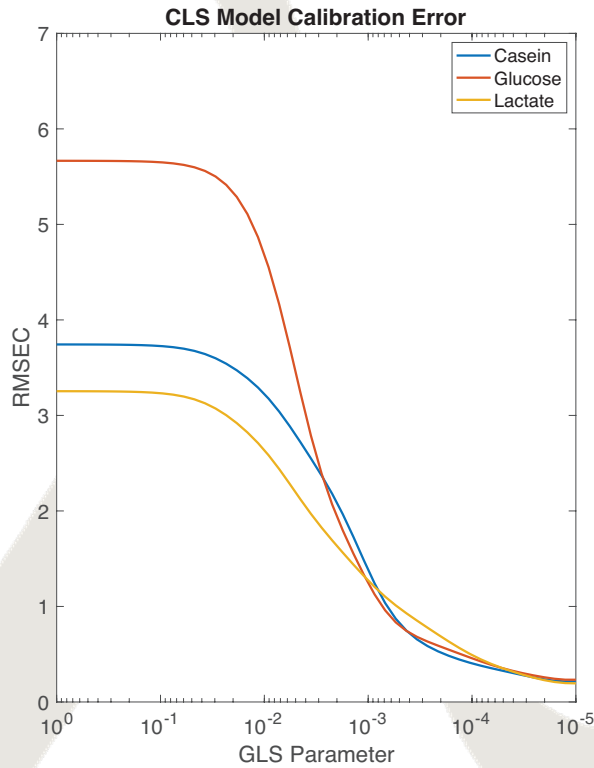- GLS generally outperforms EPO, so only GLS results shown

EIGENVECTOR
RESEARCH INCORPORATED

# NIR Data Sets



- Grain protein from Tormod Naes/Tomas Isaakson (CGL)
  - Casein, glucose, lactate and moisture, 231 samples (split 153/78), 117 wavelengths, full 3 component mixture design,

- Styrene-butadiene from Dupont/Chuck Miller (SBR)
  - Styrene, cis-, trans- and 1,2-butadiene, 70 samples (split 60/10), 141 wavelengths.

- Hydrocarbon mixture from Willem Windig (WW)
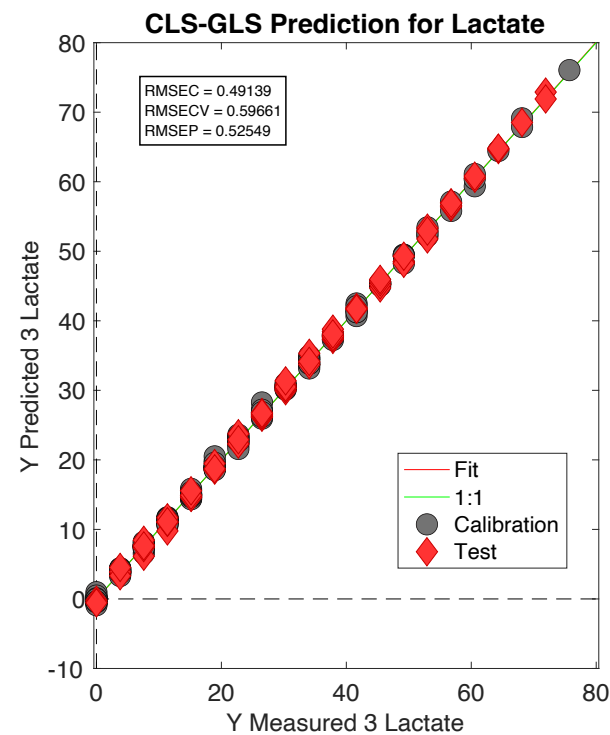  - Butanol, dichloromethane, methanol, dichloropropane and acetone, 140 samples (split 93/47), 700 wavelengths, full design.
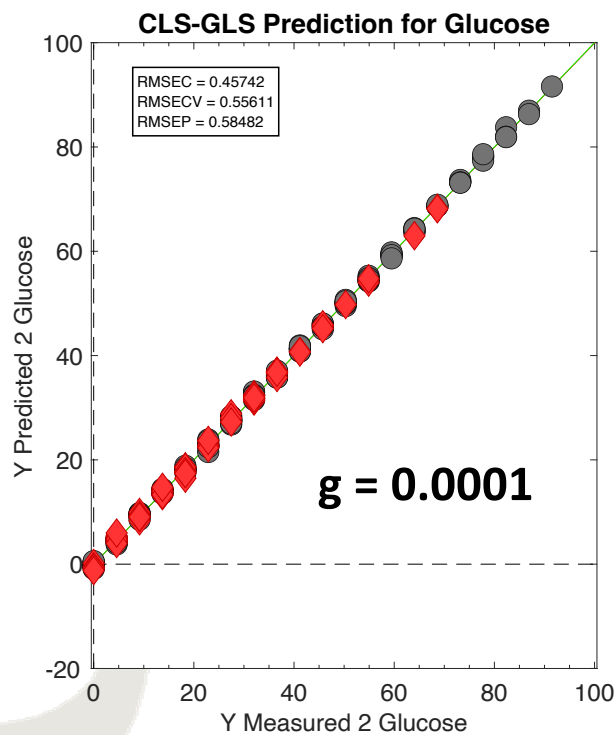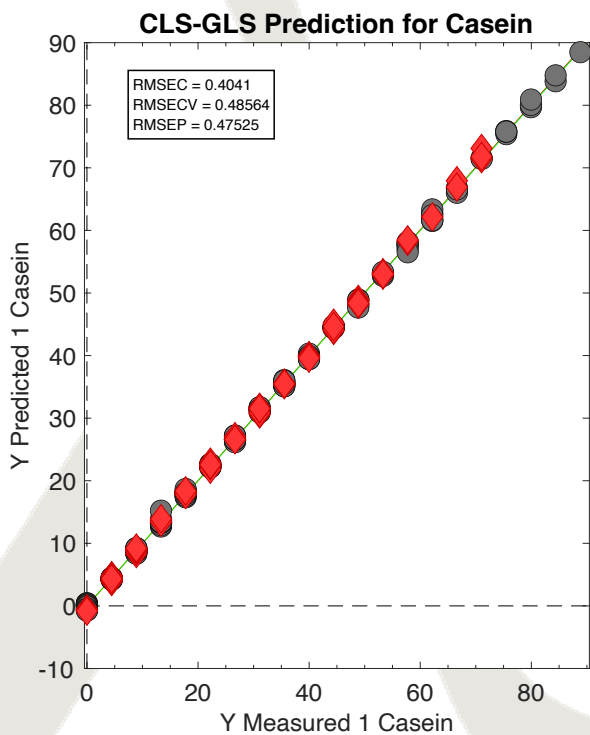
# CLS Predictions for CGL



**CLS Prediction for Casein**

RMSEC = 3.7433
RMSECV = 3.8203
RMSEP = 3.0879

**CLS Prediction for Glucose**

RMSEC = 5.6661
RMSECV = 5.82
RMSEP = 4.3923

**CLS Prediction for Lactate**

RMSEC = 3.2538
RMSECV = 3.3194
RMSEP = 3.4515

EIGENVECTOR
RESEARCH INCORPORATED

# CGL Model Error – RMSEC/CV/P

# CLS-GLS Predictions for CGL

# Diagnostic Information



**Estimated Pure Component Spectra from CLS**

Legend:
- Casein (31.78%)
- Glucose (56.01%)
- Lactate (12.21%)

**First 4 PCs from Half-Residuals**

Legend:
- PC 1 (93.75%)
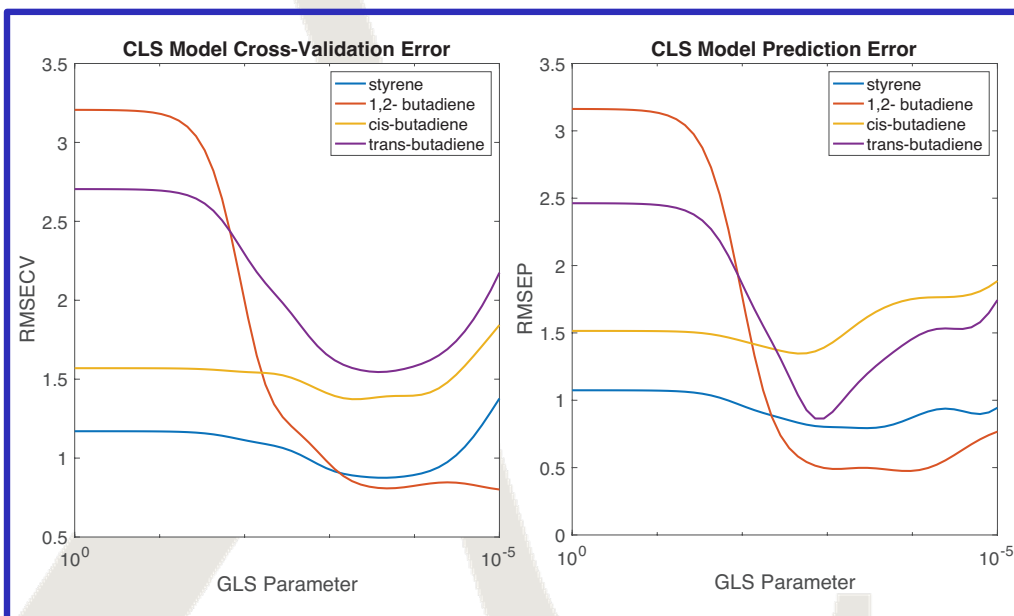- PC 2 (4.17%)
- PC 3 (0.99%)
- PC 4 (0.71%)

EIGENVECTOR RESEARCH INCORPORATED

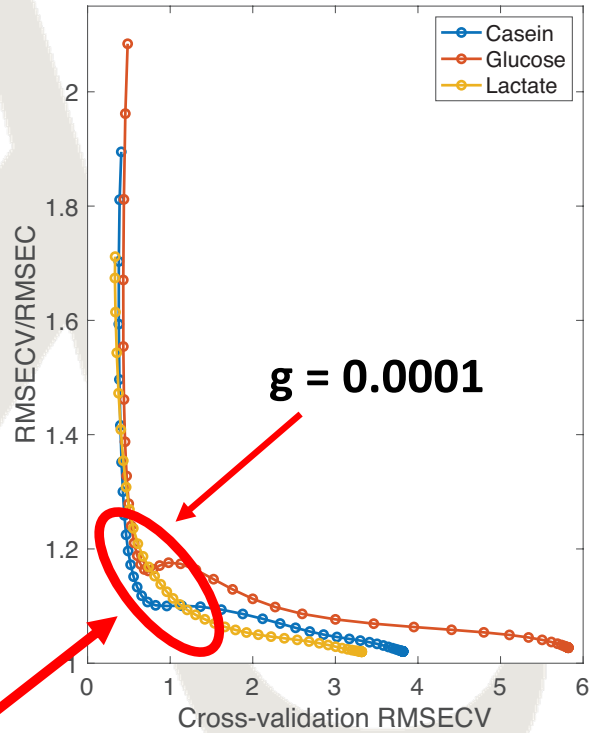# SBR Data Results

# WW Data Results
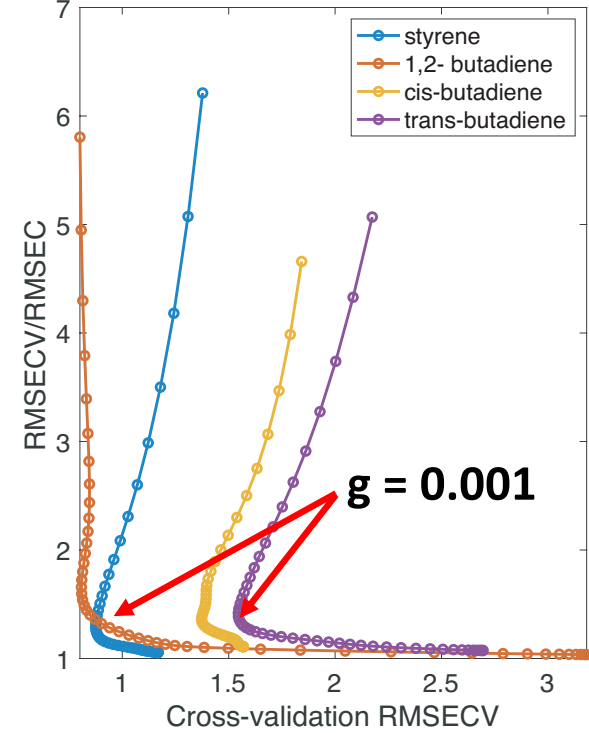
# Choosing Meta-parameters

- Usual aspects of cross-validation apply
- Watch for overfit, *i.e.* when fit error (RMSEC) is much lower than prediction error (RMSECV)
  - Plot Overfit vs. Cross-validation error
    - Ratio RMSECV/RMSEC versus RMSECV
    - In units of predicted variable

**EIGENVECTOR**
**RESEARCH INCORPORATED**

# RMSECV/C for CGL & SBR

# Conclusions

- EPO and GLSW filters can be used to improve CLS model predictive performance – **Gray Models**
  - Key is use of "half-residuals" $R_h$
  - One adjustable parameter ($k$ or $g$)
- Resulting models competitive with PLS in predictive ability
- Model selection criteria as usual
  - Overfit (RMSECV/C) vs. Prediction (RMSECV) very useful!
- Main advantages interpretability, explain-ability
- Available in PLS_Toolbox/Solo 9.3

**EIGENVECTOR**
RESEARCH INCORPORATED