

Gray Classical Least Squares

Putting Chemistry Back Into Calibration Models

Barry M. Wise¹, Donal O'Sullivan¹ and Rasmus Bro²

¹Eigenvector Research, Inc.

²University of Copenhagen



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 propose a method for developing CLS models with predictive properties competitive with ILS formulations. This is done by creating the CLS model “half-residuals” and then using these to develop pre-filters with Generalized Least Squares Weighting (GLSW) or External Parameter Orthogonalization (EPO).



Outline

- The Classical Least Squares Model
 - Ways to generate spectral residuals
- Clutter Filters
 - EPO: External Parameter Orthogonalization
 - GLSW: Generalized Least Squares Weighting
- Toy Example: 3 components + baseline
- Example Data Sets
- Results
 - Choosing the Meta-parameters
 - Diagnostics
- Conclusions

3



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



Using the CLS Model

- If S is known, the \hat{C} can be estimated from
 - $\hat{C} = XS(S^T S)^{-1}$
- 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.

5



CLS Spectral Residuals

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

*R_h needs sexier name!
Interference revealing residuals?*

6



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. ??

7

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)

8

EPO Filter

- Given a matrix \mathbf{Z} which represents extraneous variation (matrix effects, clutter), decompose \mathbf{Z} as
 - $\mathbf{Z} = \mathbf{USV}^T$ $\mathbf{Z} = ?$
- The number of filter factors k must be specified, then
 - $\mathbf{F}_{epo} = \mathbf{I} - \mathbf{V}_k \mathbf{V}_k^T$
- \mathbf{F}_{epo} is applied to \mathbf{X} before calibration (and during prediction), and removes variations in the first k dimensions represented in \mathbf{Z}

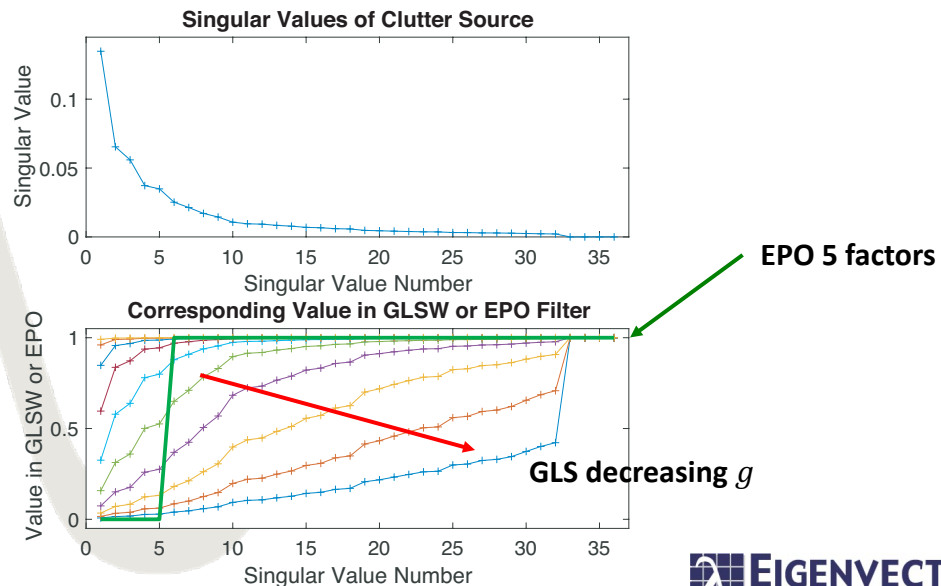
9

GLSW Filter

- Similar to EPO filter except it shrinks dimensions rather than completely eliminating them
- Starting the the decomposition of \mathbf{Z} above then
 - $\mathbf{F}_{glsw} = \mathbf{VD}^{-1}\mathbf{V}^T$
- Where the diagonal elements of \mathbf{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 \mathbf{S} and g is a tune-able parameter which controls the shrinkage

10

Comparison of EPO & GLSW



11

Combining CLS & Filters

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

CRUX

12

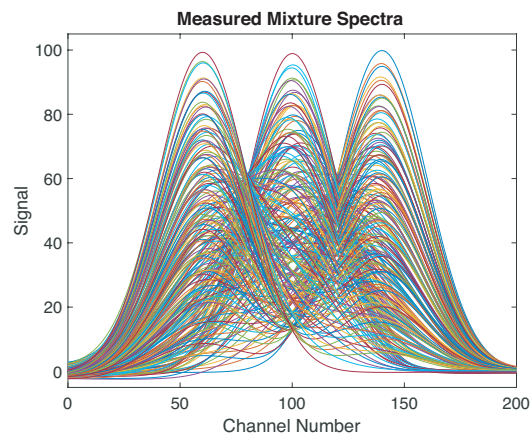
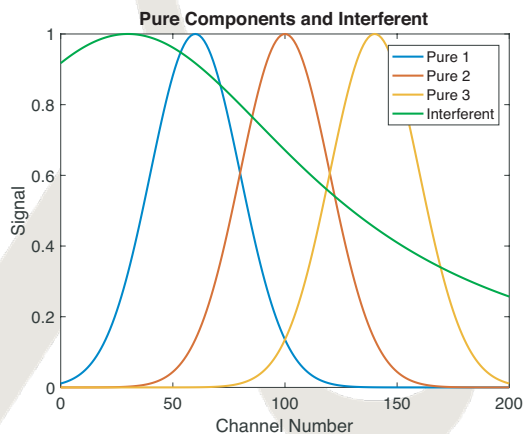
The CLS Gray Model

- 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)

13



Toy Example: 3 Components + Baseline

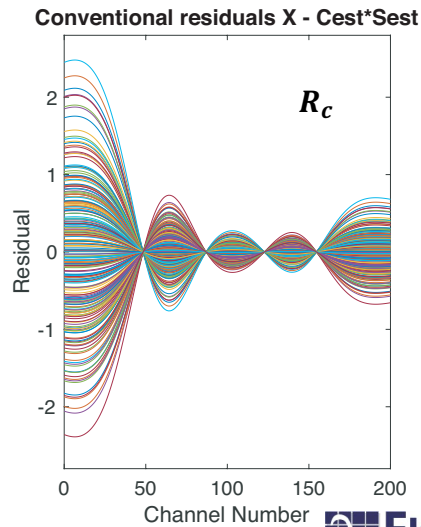
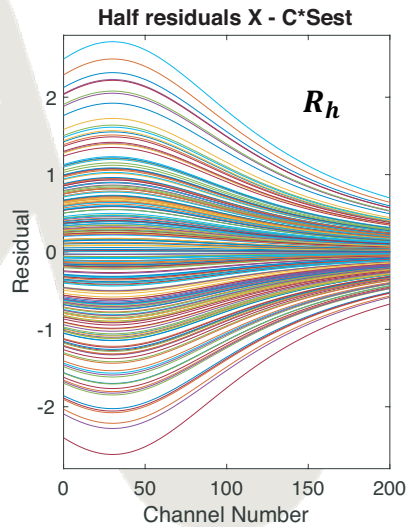


14

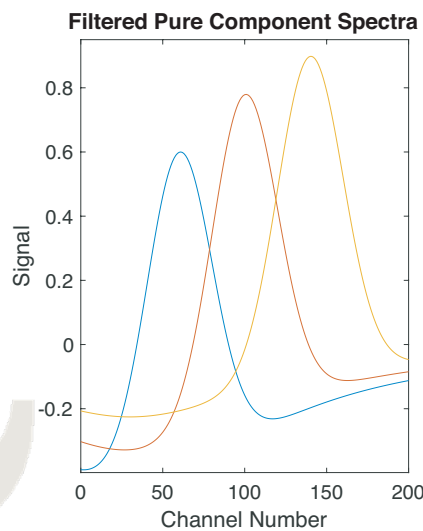
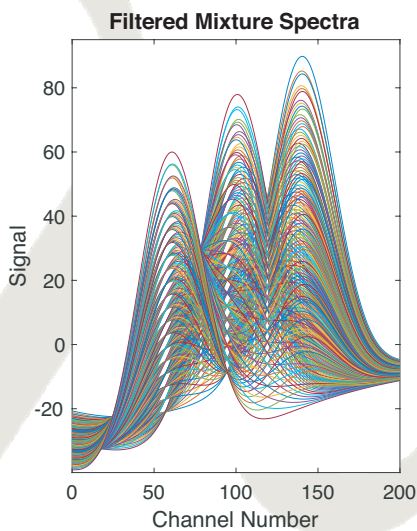


Comparison of Residuals

Unmodeled
interferent
obvious in R_h



EPO (1) Filtered Spectra

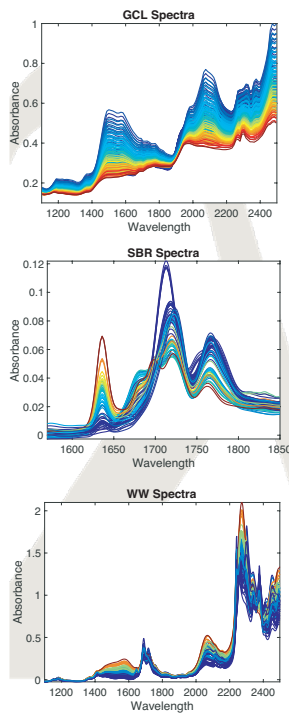


RMSEC and RMSEP
for filtered data now
at machine precision

EMM & EPO Equivalence

- In the Extended Mixture Model (EMM), clutter factors \mathbf{V}_k are augmented onto estimated spectra during prediction
 - $\mathbf{S}_a = [\mathbf{S} \mid \mathbf{V}_k]$
 - $[\mathbf{C}_n \mid \mathbf{C}_f] = \mathbf{X}_n \mathbf{S}_a (\mathbf{S}_a^T \mathbf{S}_a)^{-1}$
- Where \mathbf{C}_n are the analyte concentrations, \mathbf{C}_f are the pseudo concentrations of the clutter factors
- **Mathematically identical** to using EPO as a prefilter
- No direct equivalent for GLSW

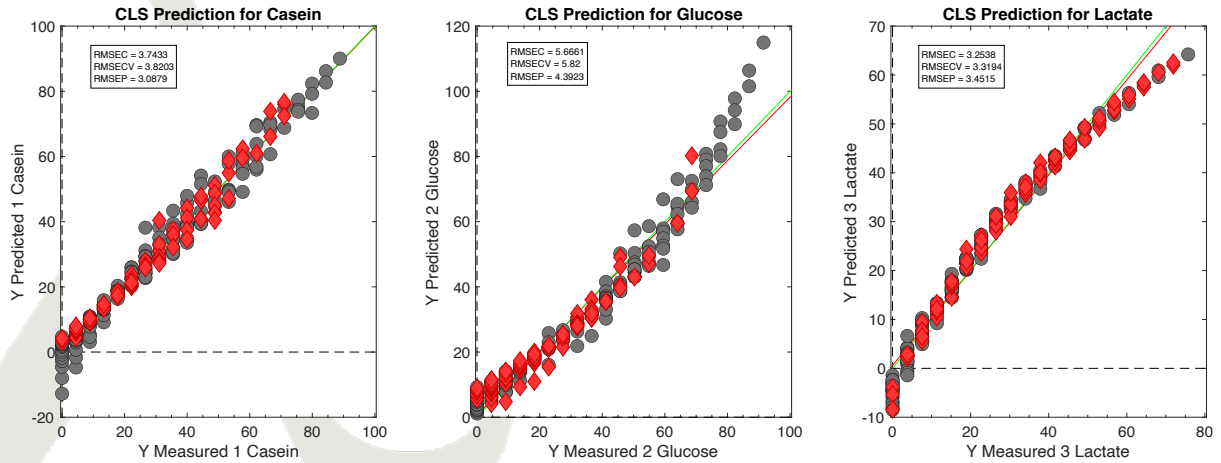
17



Example 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.

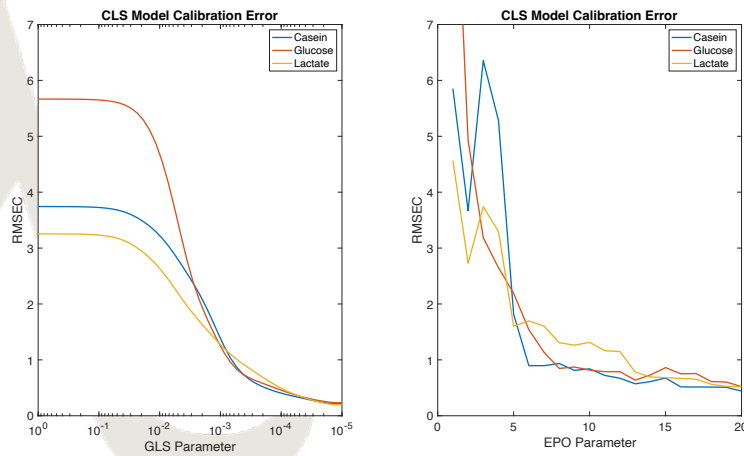
CLS Predictions for CGL



19



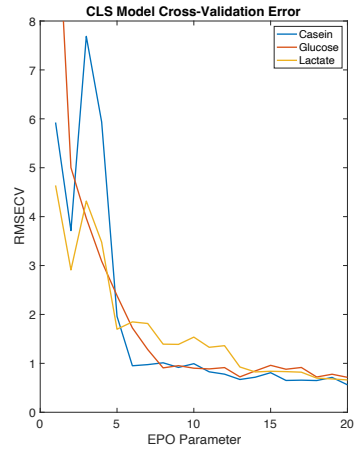
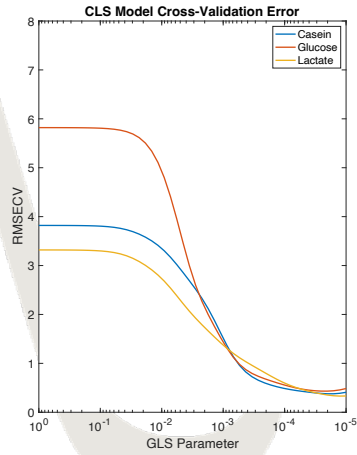
CGL Calibration Error - RMSEC



20

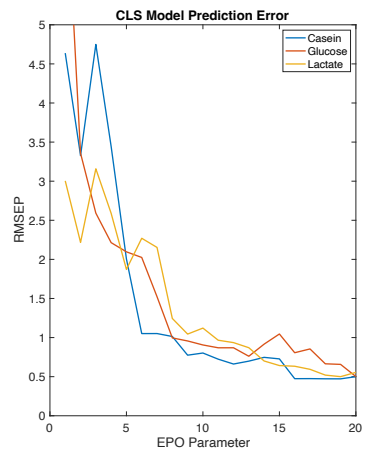
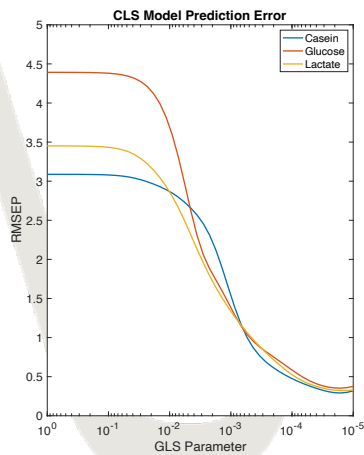


CGL Cross-validation Error - RMSECV



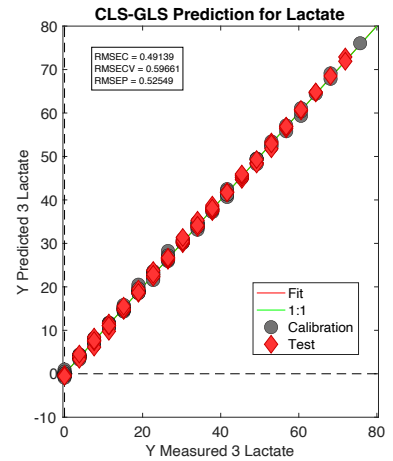
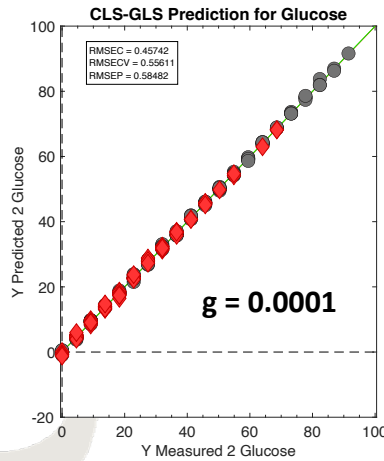
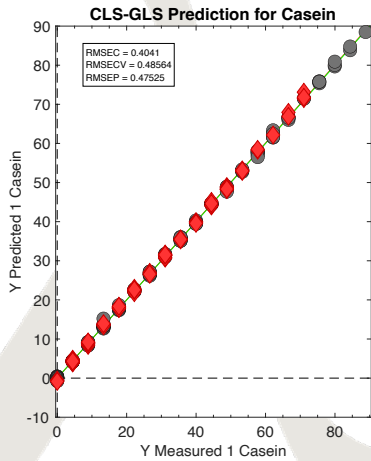
21

CGL Prediction Error - RMSEP



22

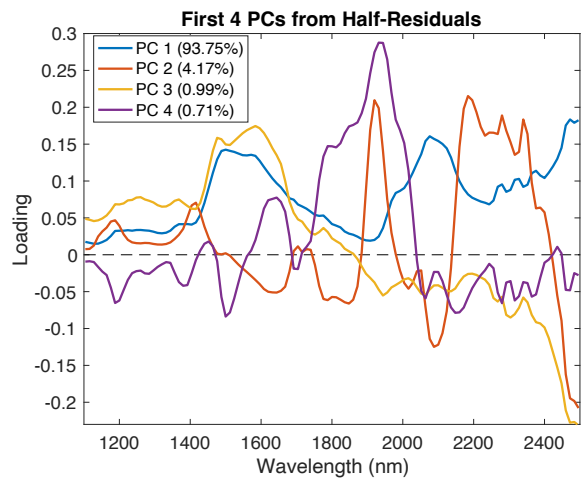
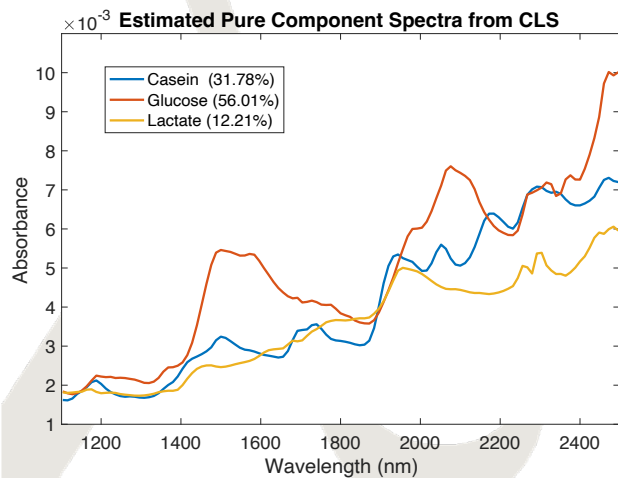
CLS-GLS Predictions for CGL



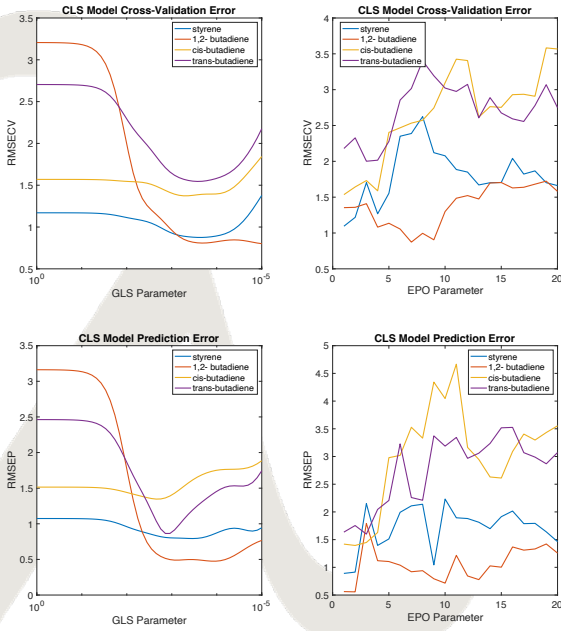
CLS 0.47 0.58 0.53 $g=0.0001$
PLS 0.52 0.68 0.59 12 LVs



Diagnostic Information

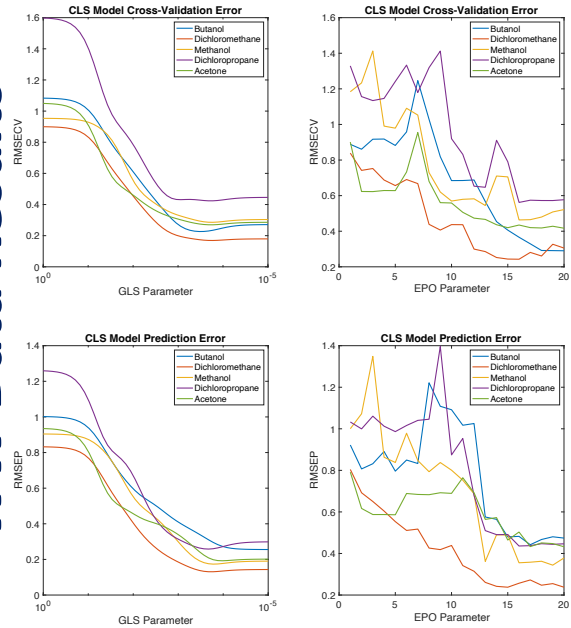


SBR Data Results



25

WW Data Results



 **EIGENVECTOR**
RESEARCH INCORPORATED

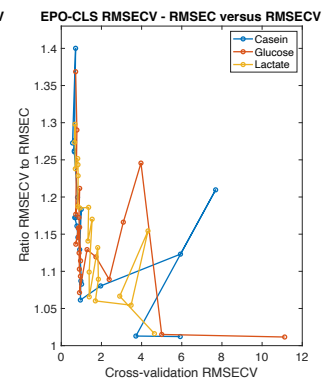
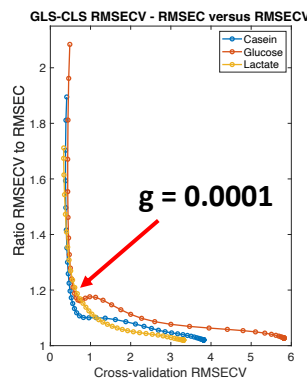
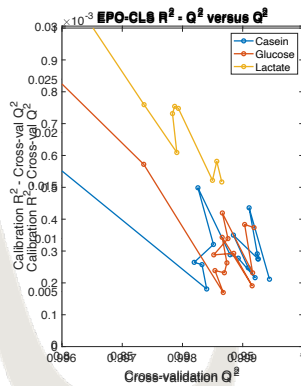
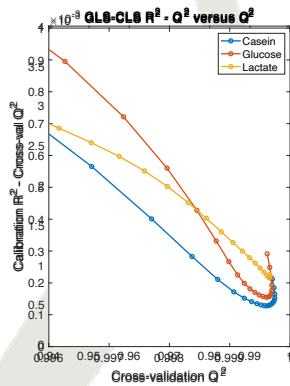
Choosing Meta-parameters

- Usual aspects of cross-validation apply
- Cross-validation appears to work as usual
- Watch for overfit (difference between fit and prediction)
- Looked at $R^2 - Q^2$ versus Q^2
 - Q^2 = Cross-validation R^2
- Also RMSECV/RMSEC versus RMSECV
 - Essentially the same information
 - In units of predicted variable
 - *Much* easier to use!

26

 **EIGENVECTOR**
RESEARCH INCORPORATED

R²/Q² and RMSEC/CV for CGL



27



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
 - Prefer RMSEC/CV instead R²/Q²
- Main advantages interpretability, explain-ability
- Available in PLS_Toolbox/Solo 9.0

28

