

# Multivariate Calibrations with Auxiliary Information

Tong Peifeng

Guanghua School of Management, Peking University

Supervised by Song Xi Chen and Cheng Yong Tang  
2023/04/16



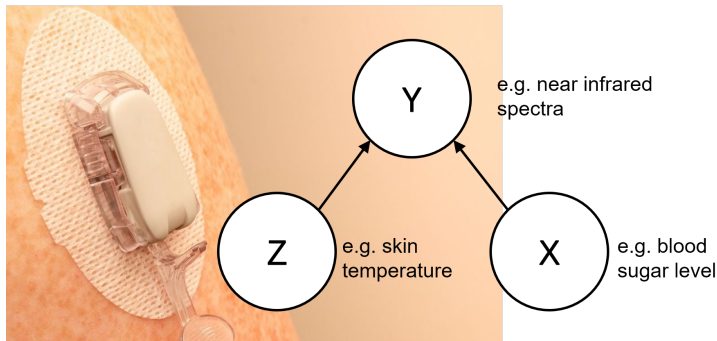
## 1 Motivation of calibration problem

## 2 Main results

## 3 Simulation results

## 4 Case studies

# Introduction of inverse regression



Inverse prediction is a common situation in many areas such as economics, sociology and earth sciences. In such a prediction setting, we desire to use some cheap and quick but **error-prone** measurement  $Y$ , to predict the true amount of a constituent  $X$ , which in itself can be measured accurately in laboratory condition but **much greater effort or cost** (Brown, 2006).

## Bias due to conditioning on a collider

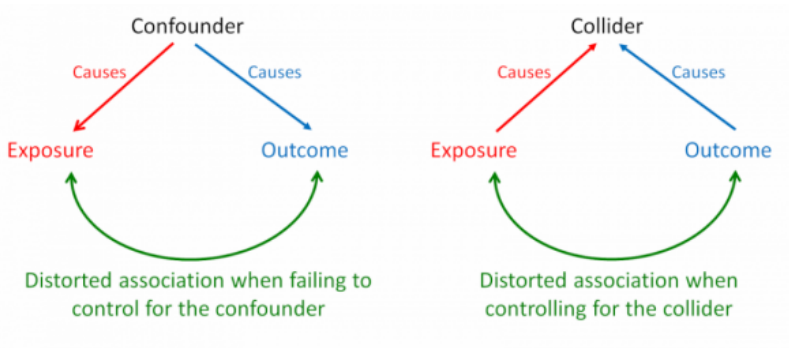


Figure 1: An illustration of the differences between confounder and collider.

## The calibration model setting

Given a sample of  $n$  observations  $\{(\mathbf{y}_i, \mathbf{x}_i, \mathbf{z}_i)\}_{i=1}^n$ , where  $\mathbf{y}_i \in \mathbb{R}^q$  is the multivariate response with covariates  $\mathbf{x}_i \in \mathbb{R}^p$  and  $\mathbf{z}_i \in \mathbb{R}^{p'}$ . We consider the following two models:

**Observation model:**

$$\mathbf{y}_i = \boldsymbol{\alpha} + \mathbf{B}^T \mathbf{x}_i + \mathbf{D}^T \mathbf{z}_i + \mathbf{e}_i, \quad (1)$$

**Calibration model:**

$$\mathbf{y}'_j = \boldsymbol{\alpha} + \mathbf{B}^T \boldsymbol{\xi} + \mathbf{D}^T \mathbf{z}'_j + \mathbf{e}'_j. \quad (2)$$

In the calibration model (2), repeated measurements  $\{\mathbf{y}'_j\}_{j=1}^l$  are available together with **only**  $\{\mathbf{z}'_j\}_{j=1}^l$ . The un-observable counterpart of  $\mathbf{x}'_j$  is considered as shared by these  $l$  measurements, and is considered as a standalone parameter denoted by  $\boldsymbol{\xi} \in \mathbb{R}^p$ .

## The calibration model setting

Collectively, we write the two models in their matrix forms as

$$Y = \mathbf{1}_n \alpha^T + XB + ZD + E, \quad (3)$$

$$Y' = \mathbf{1}_l \alpha^T + \mathbf{1}_l \xi^T B + Z'^T D + E', \quad (4)$$

where residuals  $E$  and  $E'$  are the exogenous residuals whose row vectors  $e_i^T$  are assumed to be independent and identically distributed satisfying

$$\mathbb{E}(e_i) = 0, \mathbb{E}(e_i e_i^T) = \Gamma. \quad (5)$$

We begin with this conventional cross-sectional setting and then extend the calibration problem to cover dependent  $y$ 's.

- 1 Motivation of calibration problem
- 2 Main results
- 3 Simulation results
- 4 Case studies

## Two traditional ways in solving $\xi$

The so-called generalized least squares estimator (GLS) and the inverse regression (IR) estimator are generalized to our new setting with auxiliary variables  $Z$  in the calibration problem.

**GLS:**

$$\hat{\xi}_{glS} = (\hat{B}S^{-1}\hat{B}^T)^{-1}\hat{B}S^{-1}(\bar{y}' - \hat{\alpha} - \hat{D}^T\bar{z}'), \quad (6)$$

where  $\hat{\alpha}$ ,  $\hat{B}$ , and  $\hat{D}$  are the OLS estimators of the observation model (3),  $\bar{z}' = (l^{-1}\mathbf{1}_l^T Z')^T$ ,  $S = (Y - \hat{Y})^T(Y - \hat{Y})$ , and  $\hat{Y} = \mathbf{1}_n\hat{\alpha}^T + X\hat{B} + Z\hat{D}$  is the fitted/predicted responses of the model (3).

**IR:**

$$\hat{\xi}_{ir} = \hat{\theta} + \hat{\Phi}^T\bar{y}' + \hat{\Psi}^T\bar{z}', \quad (7)$$

where  $\hat{\theta}$ ,  $\hat{\Phi}$  and  $\hat{\Psi}$  are the OLS estimators of a working model  $X = \mathbf{1}_n\theta^T + Y\Phi + Z\Psi + \tilde{E}$ .



## A Bayesian view for the IR estimator

We denote by  $\mathbb{P}_A = A(A^T A)^{-1}A^T$  and  $\mathbb{M}_A = I - \mathbb{P}_A$  be the projection matrices to its column space and the orthogonal complement, respectively.

Let  $H_{ir} = (\hat{B}S^{-1}\hat{B}^T)^{-1}(X^T\mathbb{M}_Z X)^{-1}$  and  $\zeta = X^T Z(Z^T Z)^{-1}\bar{z}'$ . Our main results on the properties of the GLS and inverse regression estimators, and a new **shrinkage estimator** are given in the following theorems.

### Theorem

(i) Under some regularity conditions,

$$\hat{\xi}_{ir} = (I + H_{ir})^{-1}(\hat{\xi}_{glS} - \zeta) + \zeta. \quad (8)$$

(ii) Under the normality assumption and a proper prior distribution,  $\hat{\xi}_{ir}$  is the **Bayesian estimator**, namely the posterior mean that minimizes the Bayesian risk.

The prior distribution of  $\xi$  given  $X$ ,  $Z$  and  $Z'$  is the multivariate t-distribution  $T_{\nu-p}(\zeta, \frac{1}{\nu-p}(l^{-1} + n^{-1} + c_3 - \frac{1}{4}c_2^T C_1^{-1} c_2)C_1^{-1})$ , where  $\nu = n - p - p' - q$ ,  $C_1 = (X^T \mathbb{M}_Z X)^{-1}$ ,  $c_2 = ((X^T \mathbb{M}_Z X)^{-1} X^T Z(Z^T Z)^{-1} + (X^T X)^{-1} X^T Z(Z^T \mathbb{M}_X Z)^{-1})\bar{z}'$  and  $c_3 = \bar{z}'^T (Z^T \mathbb{M}_X Z)^{-1} \bar{z}'$ .

## A better shrinkage estimator

We define the limiting predictive mean square error of the  $\hat{\xi}$  as

$$G(\xi, \zeta) = \|\{(I + H)^{-1} - I\}(\xi - \zeta)\|_2^2 + l^{-1} \text{tr}\{(I + H)^{-1} (B\Gamma^{-1}B^T)^{-1} (I + H^T)^{-1}\}.$$

Upon treating the design as random, we define the averaged mean squared error as  $\text{AMSE} = \mathbb{E}\{G(\xi, \zeta)\}$ .

Then, for the class of estimators  $\hat{\xi} = (I + H)^{-1}(\hat{\xi}_{\text{gls}} - \zeta) + \zeta$  with a positive definite  $H \in \mathbb{R}^{p \times p}$ , we propose to find the **optimal shrinkage estimator** that minimizes the AMSE; the result is given in the next theorem.

### Theorem

*If  $\xi$  and  $\{\mathbf{x}_i\}_{i=1}^n$  are independent and identically distributed (IID) from a distribution having finite  $(2 + \delta)$ -th moment for a  $\delta > 0$ , let the covariance matrices of  $X$  and  $Z$  be  $\Gamma_x$ ,  $\Gamma_z$  and  $\Gamma_{xz}$ , the AMSE is minimized at*

$$\tilde{H}_{\text{opt}} = l^{-1} (B\Gamma^{-1}B^T)^{-1} (\Gamma_x - \Gamma_{xz}\Gamma_z^{-1}\Gamma_{xz}^T)^{-1}.$$

*Furthermore, as  $n \rightarrow \infty$ ,*

$$H_{\text{opt}} = l^{-1} (\hat{B}\hat{S}^{-1}\hat{B}^T)^{-1} (X^T \mathbb{M}_Z X)^{-1}$$

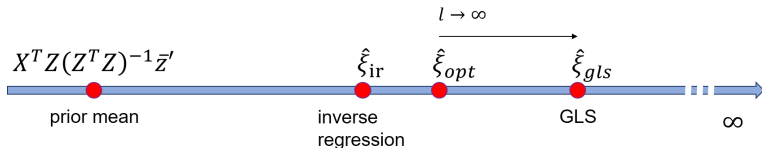
*is consistent to  $\tilde{H}_{\text{opt}}$ .*

## The relationship between the proposed estimators

Here  $\tilde{\mathbf{H}}_{opt}$  in Theorem 2 is an unknown hyper-parameter; and we propose to implement its empirical counterpart  $\mathbf{H}_{opt}$  for empirically obtaining the optimal shrinkage estimator:

$$\hat{\xi}_{opt} = (\mathbf{I} + \mathbf{H}_{opt})^{-1}(\hat{\xi}_{gls} - \zeta) + \zeta. \quad (9)$$

In summary, the relationship between these estimators is depicted in Figure 2, which shows the spectrum of trade-offs between bias and variance in the class of shrinkage estimators.



**Figure 2:** The relationship among GLS, inverse regression and optimal shrinkage estimators. When  $l = 1$ , the optimal shrinkage estimator coincides with the inverse regression estimator. The optimal shrinkage estimator converges to the GLS estimator as  $l \rightarrow \infty$ .

## An objective function interpretation

Remarkably, as shown in the next theorem, the optimal shrinkage estimator, the GLS and inverse regression estimators link to a class of the **penalized estimations**:

$$\min_{\xi} \text{tr}[\{Y' - \hat{Y}'(\xi)\}S^{-1}\{Y' - \hat{Y}'(\xi)\}^T] + \lambda\sigma^2(\xi), \quad (10)$$

where  $\hat{Y}'(\xi) = \mathbf{1}_l\hat{\alpha}^T + \mathbf{1}_l\xi^T\hat{B} + Z'\hat{D}$ ,  $\sigma^2(\xi) = 1/l + 1/n + \xi^T C_1 \xi - c_2^T \xi + c_3$ , with  $C_1$ ,  $c_2$  and  $c_3$  defined in Assumption 2 (b).

### Theorem

*Under some regularity conditions,*

- ①  $\hat{\xi}_{ir}$  is the solution for  $\lambda = l$ ;
- ②  $\hat{\xi}_{opt}$  is the solution for  $\lambda = 1$ ;
- ③  $\hat{\xi}_{glS}$  is the solution for  $\lambda = 0$ .

## Calibration with random effect

So far, our method treats the  $n$  measurements in  $Y$  of model (3) as independent. For incorporating the **correlated scenarios**, we extend the framework to accommodate the mixed or random effects in handling broader situations. Specifically, we consider an extended observation model

$$Y = \mathbf{1}_n \alpha^T + XB + ZD + \Psi A + E, \quad (11)$$

where  $A = (\mathbf{a}_1, \dots, \mathbf{a}_k)^T \in \mathbb{R}^{k \times q}$ ,  $\mathbf{a}_j \in \mathbb{R}^q$  ( $j = 1, \dots, k$ ) are  $k$  random effects, and  $\Psi \in \mathbb{R}^{n \times k}$  is a **known model matrix**. We assume  $\mathbf{a}_k \sim \mathcal{N}(0, \Gamma_1)$  and  $\mathbf{e}_i \sim \mathcal{N}(0, \Gamma_2)$ . We develop an expectation-maximization (EM) algorithm that is computationally highly efficient for the estimation of  $\alpha$ ,  $B$ ,  $D$ ,  $\Gamma_1$  and  $\Gamma_2$ . The vector form of (11) is

$$\text{vec}(Y^T) = \text{vec}(\alpha \mathbf{1}_n^T) + \text{vec}(B^T X^T) + \text{vec}(D^T Z^T) + (\Psi \otimes I) \text{vec}(A^T) + \text{vec}(E^T). \quad (12)$$

Model (11) implies that

$$\text{var}\{\text{vec}(Y^T)\} = \Psi \Psi^T \otimes \Gamma_1 + I_n \otimes \Gamma_2.$$

## Calibration with random effect

Corresponding to (2), the calibration model with random effect is

$$y'_j = \alpha + B^T \xi + D^T z'_j + A^T \psi' + e'_j \quad \text{for } j = 1, 2, \dots, l, \quad (13)$$

where the extra  $A^T \psi'$  reflects the reality that some of the random effects may be shared as reflected by the design  $\psi'$ .

Let  $\hat{\alpha}, \hat{B}, \hat{D}, \hat{\Gamma}_1$ , and  $\hat{\Gamma}_2$  be the estimators using the EM algorithm. Then, the **GLS estimator** incorporating the random effects is

$$\hat{\xi}_{gl\bar{s}} = (\hat{B}S^{-1}\hat{B}^T)^{-1}\hat{B}S^{-1}(\bar{y}' - \hat{\alpha} - \hat{D}^T\bar{z}' - \hat{A}^T\psi'), \quad (14)$$

where  $\hat{A} = \mathbb{E}(A^T|Y, \hat{\alpha}, \hat{B}, \hat{D}, \hat{\Gamma}_1, \hat{\Gamma}_2)$  and

$$S = l^{-1}\hat{\Gamma}_2 + \psi'^T\psi' \otimes \hat{\Gamma}_1 - (\psi'^T\Psi^T \otimes \hat{\Gamma}_1)(\Psi\Psi^T \otimes \hat{\Gamma}_1 + I_n \otimes \hat{\Gamma}_2)^{-1}(\Psi\psi' \otimes \hat{\Gamma}_1)$$

is the conditional variance of  $\bar{y}' - \hat{\alpha} - \hat{D}^T\bar{z}' - \hat{A}\psi'$  given  $(Y, \hat{\alpha}, \hat{B}, \hat{D}, \hat{\Gamma}_1, \hat{\Gamma}_2)$ .

Inspired by (9), we can develop an **optimal shrinkage estimator** as

$$\hat{\xi}_{opt} = (I + H)^{-1}\hat{\xi}_{gl\bar{s}}, \quad (15)$$

where  $H = (\hat{B}S^{-1}\hat{B}^T)^{-1}(\frac{1}{n}X^TX)^{-1}$ .

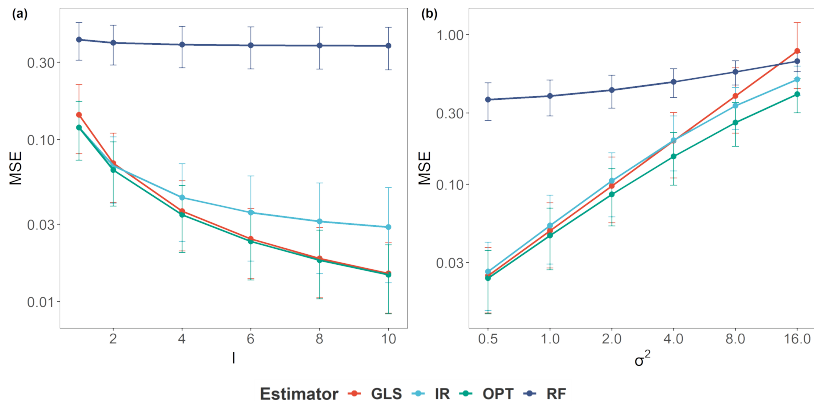
1 Motivation of calibration problem

2 Main results

3 Simulation results

4 Case studies

## A baseline setting



**Figure 3:** Empirical means and the 10th, 90th quantiles of the MSEs (in log-scale) calculated from 1000 simulation replications for the GLS, the inverse regression (IR), the optimal shrinkage (OPT) and the inverse regression with the random forest (RF) estimators with respect to numbers of repeated measurements  $l$  (a) and the variances of the noises  $\sigma^2$  (b).



# Sensitivity analysis

**Table 1:** Empirical means and standard deviations of the MSEs from 1000 simulations for the GLS, the inverse regression (IR), the optimal shrinkage (OPT) and the inverse regression with the random forest (RF) estimators with respect to (a) the number of dependence variables  $q^*$  used for training, (b) the number of covariates  $p'^*$  used for training,

(a) with respect to number of dependence variables

$q^*$	$0.2q$	$0.4q$	$0.6q$	$0.8q$	$q$
GLS	$3922 \pm 60948$	$0.298 \pm 0.568$	$0.117 \pm 0.207$	$0.069 \pm 0.045$	$0.049 \pm 0.022$
IR	$0.438 \pm 0.166$	$0.199 \pm 0.112$	$0.112 \pm 0.066$	$0.074 \pm 0.040$	$0.053 \pm 0.024$
OPT	$0.385 \pm 0.165$	$0.159 \pm 0.094$	$0.090 \pm 0.052$	$0.061 \pm 0.030$	$0.045 \pm 0.018$
RF	$0.746 \pm 0.120$	$0.599 \pm 0.114$	$0.489 \pm 0.103$	$0.432 \pm 0.094$	$0.391 \pm 0.085$

(b) with respect to number of covariates

$p'^*$	$0.2p'$	$0.4p'$	$0.6p'$	$0.8p'$	$p'$
GLS	$0.315 \pm 0.255$	$0.187 \pm 0.171$	$0.110 \pm 0.086$	$0.070 \pm 0.044$	$0.049 \pm 0.024$
IR	$0.211 \pm 0.085$	$0.149 \pm 0.069$	$0.103 \pm 0.052$	$0.072 \pm 0.036$	$0.054 \pm 0.025$
OPT	$0.214 \pm 0.097$	$0.141 \pm 0.072$	$0.092 \pm 0.049$	$0.063 \pm 0.030$	$0.046 \pm 0.019$
RF	$0.413 \pm 0.092$	$0.410 \pm 0.089$	$0.389 \pm 0.087$	$0.388 \pm 0.084$	$0.389 \pm 0.083$

## Sensitivity analysis

**Table 2:** Empirical means and standard deviations of the MSEs from 1000 simulations for the GLS, the inverse regression (IR), the optimal shrinkage (OPT) and the inverse regression with the random forest (RF) estimators with respect to (c) the sample size  $N$  and (d) the explanation power of  $X$  as represented by shrinking  $B^*$  to  $\gamma B^*$ .

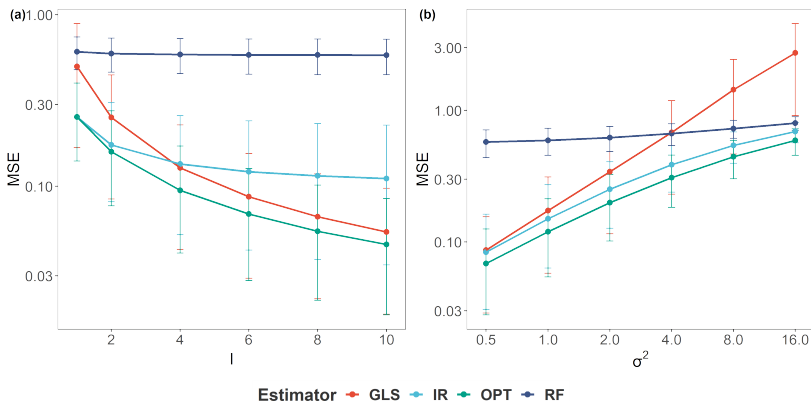
(c) with respect to sample size

$N$	30	100	300	1000	2500
GLS	$0.461 \pm 0.314$	$0.079 \pm 0.038$	$0.056 \pm 0.026$	$0.050 \pm 0.024$	$0.049 \pm 0.023$
IR	$0.418 \pm 0.245$	$0.079 \pm 0.033$	$0.059 \pm 0.026$	$0.054 \pm 0.025$	$0.053 \pm 0.024$
OPT	$0.441 \pm 0.279$	$0.074 \pm 0.030$	$0.052 \pm 0.021$	$0.047 \pm 0.019$	$0.045 \pm 0.018$
RF	$0.843 \pm 0.163$	$0.666 \pm 0.099$	$0.545 \pm 0.086$	$0.447 \pm 0.083$	$0.390 \pm 0.081$

(d) with respect to explaining power of  $X$

$\gamma$	0.1	0.2	0.33	0.5	1
GLS	$4.471 \pm 1.832$	$1.157 \pm 0.471$	$0.420 \pm 0.171$	$0.187 \pm 0.076$	$0.047 \pm 0.019$
IR	$0.863 \pm 0.040$	$0.608 \pm 0.076$	$0.356 \pm 0.080$	$0.193 \pm 0.061$	$0.051 \pm 0.021$
OPT	$0.792 \pm 0.054$	$0.498 \pm 0.078$	$0.276 \pm 0.067$	$0.150 \pm 0.046$	$0.044 \pm 0.016$
RF	$0.999 \pm 0.022$	$0.934 \pm 0.034$	$0.819 \pm 0.058$	$0.676 \pm 0.077$	$0.383 \pm 0.082$

## Spatial-temporal settings



**Figure 4:** Empirical means and the 10th, 90th quantiles of the MSEs (in the log-scale), in the presence of the random effects, calculated from 1000 simulation replications for the GLS, the inverse regression (IR), the optimal shrinkage (OPT) and the inverse regression with the random forest (RF) estimators with respect to the numbers of repeated measurements  $l$  (a) and different scales of noises  $\sigma^2$ .

- 1 Motivation of calibration problem
- 2 Main results
- 3 Simulation results
- 4 Case studies**

## The aqueous glucose data

Diabetes Mellitus (DM) is one of the common chronic diseases throughout the globe. A continuous blood glucose monitoring is important for DM patients to understand their disease progression and adjust their lifestyle management. However, such blood glucose measurements are often invasive and expensive, thus not proper for the regular use. Nowadays, near infrared spectroscopy (NIR) is a promising technique for continuous blood glucose monitoring, since it is cheap, non-invasive and easy to deploy on some wearable devices.



## The aqueous glucose data

In this case study, we applied our calibration method in a open sourced NIR spectroscopy data of aqueous glucose (Fuglerud, 2021). This dataset contained 127 samples measured under laboratory conditions. Five covariates (lactates, ethanol, caffeine, acetaminophen and temperature) were introduced to the dextrose water, alone with 4200 NIR wavelength channels ranged from 400nm-2500nm. The NIR spectral were collected in triplicate in each setting, sampled every 0.5nm frequency resolution with a bandwidth of 8.75nm.

According to Fuglerud et al. (2021), we removed those NIR bands which were highly absorbing water peaks (1900-2100nm, 2300-2500nm), the overlap of the two spectrometer detectors (1090-1110nm), and the fringes of the detectors ( $<500$  and  $>2300$  nm). After that, we selected the remaining NIR bands every 25nm, made the dimension of  $Y$  to be  $127 \times 60$ . All the data were pre-treated to have mean 0 and variance 1.

## The aqueous glucose data

**Table 3:** MSE and Bias for glucose prediction under different calibration methods

	MSE		Bias	
	$l=1$	$l=3$	$l=1$	$l=3$
GLS	0.0912	0.0683	0.2038	0.1563
inverse regression	0.0891	0.0714	0.2008	0.1604
optimal shrinkage	0.0891	0.0676	0.2008	0.1555
RF	0.4411	0.4690	0.5511	0.5668

- Brown, P. J. (2006). Inverse prediction. *Encyclopedia of Environmetrics*.
- Fuglerud, S. S. (2021). Aqueous glucose measured by NIR spectroscopy.
- Fuglerud, S. S., Ellingsen, R., Aksnes, A., and Hjelme, D. R. (2021). Investigation of the effect of clinically relevant interferences on glucose monitoring using near-infrared spectroscopy. *Journal of Biophotonics*, 14(5).



*Thanks!*