

# BSTT562 Final Project I

*Ruizhe Chen, Hesen Li*

*December 14th, 2018*

## 1 Abstract

There has been a long existing problem in analytical chemistry: handling uncertainty in chemical measurements. Major challenges in this type of problem are heteroscedasticity, which results from increasing variances with increasing concentrations and inter-laboratories variations. Acknowledging the fact that uncertainty in measurements should be incorporated in computing bounds on the underlying true concentration, Bhaumik and Gibbons [2005] proposed a random-effect calibration curve model. They used data separation technique in computing the point estimates and confidence intervals for true unknown concentrations  $X$  that varies from low to high levels from multiple laboratories. In this final project, we want to reproduce their analysis procedure and results (presented in Table 4 Bhaumik and Gibbons [2005]) for a real data set used in Bhaumik and Gibbons [2005] using their proposed methods.

## 2 Data

To illustrate the properties of the point estimate and confidence region for the true concentration  $X$ , Bhaumik and Gibbons [2005] analyzed experimental data for cadmium from an interlaboratory study conducted by the Ford Motor Company (J. Phillips, personal communication). These data were generated as part of a blind interlaboratory study of laboratories that hold Michigan State Drinking Water Certifications for the parameters tested. The data has the characteristic of heteroscedasticity such that measurements at higher concentration levels have greater variability than those at lower levels.

The samples were prepared by an independent source, randomized, and submitted on a weekly basis over a 5-week period. Cadmium was analyzed by inductively coupled plasma atomic emissions spectroscopy using EPA method 200.7. The data set used for this example comprised five replicates at each of three concentrations (0, 20 and 100  $\mu\text{g}/\text{L}$ ) in each of the  $q = 5$  laboratories. Using the first replicate from the first three laboratories as the new measurement (i.e.,  $q' = 3$ ), we would like to reproduce the results of point estimates, variances, confidence intervals and simulated confidence levels for the true concentrations in each of the three cases. The data is displayed as follows:

Table 1: Interlaboratory Data for Cadmium (ug / L)

Lab	Replication	0	20	100
1	1	-3.000	10.00	92.00

Lab	Replication	0	20	100
1	2	4.000	20.00	100.00
1	3	-4.000	17.20	97.80
1	4	3.000	24.00	100.00
1	5	3.100	19.10	109.00
2	1	-0.060	17.82	90.45
2	2	0.010	17.30	87.61
2	3	0.115	16.57	85.55
2	4	-0.055	17.36	89.92
2	5	0.340	18.12	90.07
3	1	-7.400	27.10	107.40
3	2	-2.100	19.40	108.10
3	3	-11.400	9.00	83.80
3	4	-11.100	10.50	81.90
3	5	-1.400	19.30	94.20
4	1	1.000	21.00	96.00
4	2	-2.126	16.05	90.65
4	3	0.523	16.08	89.39
4	4	-2.000	17.00	91.00
4	5	-0.551	15.49	85.87
5	1	0.000	18.00	91.00
5	2	0.000	19.00	101.00
5	3	0.000	19.00	102.00
5	4	-1.000	18.70	92.70
5	5	0.038	19.79	99.88

### 3 Model

To measure the true concentration of an analyte,  $x$ , a simple and straightforward method is to propose a linear calibration model as:

$$y = \alpha + \beta x + e \quad (1)$$

with normal assumption on errors. However, because our data set has heteroscedastic variances, this simple linear models fails to explain the increasing measurement variation with increasing analyte concentration, which is also commonly observed in analytic data. On the other hand, a log-linear model, for example,  $y = xe^\eta$ , where  $\eta$  is a normal random variable with mean 0 and standard deviation  $\sigma_\eta$  also fails to explain the near-constant measurement variation of  $y$  for low true concentration level ( $x$ ) (Rocke and Lorenzato [1995]). To fix these two problems at the same time, Bhaumik and Gibbons [2005] proposed a log-normal model that combines both types of errors:

$$y_{ijk} = \alpha_i + \beta_i x_j e^{\eta_{ijk}} + e_{ijk}, \quad (2)$$

where  $y_{ijk}$  is the  $k$ -th measurement at the  $j$ -th concentration level in the  $i$ -th laboratory,  $i = 1, 2, \dots, q, j = 1, 2, \dots, r, k = 1, 2, \dots, N_{ij}$ ,  $x_j$  is the true concentration at the  $j$ -th level, and  $\alpha_i$  and  $\beta_i$  are the calibration parameters for the  $i$ -th laboratory.

In this model the  $\eta_{ijk}$ 's represent proportional error at higher true concentrations, and we assume that the distributions of the proportional errors remain the same for all laboratories.

Note that for this particular model, when  $x$  is 0 or near 0, the model is reduced to:

$$\begin{aligned} y_{ijk} &= \alpha_i + \beta_j \times 0 \times e^{\eta_{ijk}} + e_{ijk} \\ y_{ijk} &= \alpha_i + e_{ijk} \end{aligned} \quad (3)$$

## 4 Analysis

We choose the Method of Moments to estimate the model parameters since it is straightforward and easy to implement when observations with lower concentrations are available, which happens to be the case with our Interlaboratory Data for Cadmium. Besides, the estimates obtained by Method of Moments are asymptotically efficient.

### 4.1 The Parameter Estimation Procedures (Using Data Separation Technique)

#### 4.1.1 (Partial) Estimation Using Low-Concentration Observations

We estimate the variance  $\sigma_e^2$  of the additive errors  $e_{ijk}$ 's (of the  $k$ -th measurement at the  $j$ -th concentration level in the  $i$ -th laboratory,  $i = 1, 2, \dots, q, j = 1, 2, \dots, r, k = 1, 2, \dots, N_{ij}$ ) that are present primarily at low-level concentrations and the calibration parameter  $\alpha_i$  for the  $i$ -th laboratory using low-concentration observations.

Assuming that observations corresponding to zero or near-zero are available. Let  $y_{i0k}$  be the  $k$ th measured concentration corresponding to the true low-level concentration from laboratory  $i$  and  $n_{i0}$  is the number of samples with true low-level concentrations submitted to laboratory  $i$ :

- (1) Estimate  $\sigma_e^2$  by the variance of the observations with zero or near-zero concentrations ( $y_{i0k}$ ) by

$$\hat{\sigma}_e^2 = \frac{1}{q} \sum_{i=1}^q \left( \frac{\sum_{k=1}^{n_{i0}} (y_{i0k} - \bar{y}_{i0})^2}{n_{i0} - 1} \right) \quad (4)$$

(2) Estimate  $\alpha_i$  from the observations with zero or near-zero concentrations ( $y_{i0k}$ ) by

$$\hat{\alpha}_i = \frac{\sum_{k=1}^{n_{i0}} y_{i0k}}{n_{i0}} \quad (5)$$

#### 4.1.2 (Partial) Estimation Using Higher-Concentration Observations

(1) Let  $z_{ijk} = \frac{(y_{ijk} - \alpha_i)}{x_j}$  and  $u_{ijk} = \frac{e_{ijk}}{x_j}$ .

Using the estimate  $\hat{\alpha}_i$  of  $\alpha_i$  we compute the  $z_{ijk}$ 's by

$$\begin{aligned} y_{ijk} &= \alpha_i + \beta_i x_j e^{\eta_{ijk}} + e_{ijk} \\ y_{ijk} - \alpha_i &= \beta_i x_j e^{\eta_{ijk}} + e_{ijk} \\ \frac{y_{ijk} - \alpha_i}{x_j} &= \beta_i e^{\eta_{ijk}} + \frac{e_{ijk}}{x_j} \text{ (given } x_j' \text{'s have higher concentrations)} \\ z_{ijk} &= \beta_i e^{\eta_{ijk}} + \mu_{ijk} \\ \frac{z_{ijk}}{\beta_i} &= e^{\eta_{ijk}} + \frac{\mu_{ijk}}{\beta_i} \text{ (given } \beta_i \neq 0). \end{aligned} \quad (6)$$

(2) Let

$$\gamma = E(e^\eta) = e^{\sigma_\eta^2/2}, \quad \mu_{zi} = E(z_{ijk}), \quad \sigma_{zi}^2 = \frac{\sum_{j=1}^r \text{var}(z_{ijk})}{r}, \quad \sigma_\mu^2 = \frac{\sigma_e^2 \sum_{j=1}^r 1/x_j^2}{r} \quad (7)$$

Using equation (6) we can calculate

$$\frac{\mu_{zi}}{\beta_i} = \gamma \text{ and } \frac{\sigma_{zi}^2}{\beta_i^2} = (\gamma^4 - \gamma^2) + \frac{\sigma_\mu^2}{\beta_i^2}. \quad (8)$$

(3) Replacing  $\gamma$  by  $\frac{\mu_{zi}}{\beta_i}$  in the second part of equation (8) we obtain the following equation:

$$\beta_i = \sqrt{\frac{\mu_{zi}^4}{\sigma_{zi}^2 - \sigma_\mu^2 + \mu_{zi}^2}} \text{ and } \sigma_\eta^2 = \frac{2 \sum_{i=1}^q \ln(\mu_{zi}/\beta_i)}{q}. \quad (9)$$

To estimate  $\beta_i$  and  $\sigma_\eta$  from equation (9), we need estimates of  $\mu_{zi}$ ,  $\sigma_{zi}^2$ , and  $\sigma_\mu^2$  (the error variance of the model (6)).

(4) Let  $n_{ij}$  be the number of observations with higher concentrations collected from the  $i$ -th laboratory for concentration level  $j$  and  $n_i = \sum_{j=1}^r n_{ij}$ , then

$$\begin{aligned}\hat{\mu}_{z_{ij}} &= \frac{\sum_{k=1}^{n_{ij}} z_{ijk}}{n_{ij}}, \quad \hat{\mu}_{z_i} = \frac{\sum_{j=1}^r \hat{\mu}_{z_{ij}}}{r}, \\ \hat{\sigma}_{z_i}^2 &= \frac{\sum_{j=1}^r \sum_{k=1}^{n_{ij}} (z_{ijk} - \hat{\mu}_{z_{ij}})^2 / (n_{ij} - 1)}{r}, \quad \hat{\sigma}_{\mu}^2 = \frac{\sum_{j=1}^r \hat{\sigma}_e^2 / x_j^2}{r}\end{aligned}\tag{10}$$

- (5) Therefore We can calculate the estimates of the variance of the proportional error at higher concentrations  $\sigma_{\eta}^2$  and the calibration parameter  $\beta_i$  for the  $i$ -th laboratory as

$$\hat{\beta}_i = \sqrt{\frac{\hat{\mu}_{z_i}^4}{\hat{\sigma}_{z_i}^2 - \hat{\sigma}_{\mu}^2 + \hat{\mu}_{z_i}^2}} \quad \text{and} \quad \hat{\sigma}_{\eta}^2 = \frac{2 \sum_{i=1}^q \ln(\hat{\mu}_{z_i} / \hat{\beta}_i)}{q}\tag{11}$$

## 4.2 Point Estimation of $X$

Suppose that the same analyte with an unknown true concentration  $X$  is tested in  $q'$  independent laboratories, and that  $Y_1, Y_2, \dots, Y_{q'}$  are the corresponding new observations. We want to compute a point estimate of  $X$  using the available information from each of the  $q'$  laboratories and then combine them. Following model (2), we can obtain

$$E(Y_i) = \alpha_i + \beta_i X \gamma,$$

and hence an estimate of  $X$  from the  $i$ -th laboratory, denoted by  $\hat{X}_i$ , is

$$\hat{X}_i = \frac{Y_i - \hat{\alpha}_i}{\hat{\beta}_i \hat{\gamma}} \quad \text{and} \quad \hat{X} = \sum_{i=1}^{q'} \frac{\hat{X}_i}{q'},\tag{12}$$

where  $\hat{X}$  is the combined estimate of  $X$ .  $\hat{X}_i$  is asymptotically unbiased, and the asymptotic expression of the variance of  $\hat{X}$  is

$$\begin{aligned}var(\hat{X}) &= E[var(\hat{X}/y)] + var[E[\hat{X}/y]] \\ &= \frac{\sum_{i=1}^{q'} \sigma_e^2 / (\beta_i^2 \gamma^2) (1 + 1/n_{0i})}{q'^2} + \frac{X^2 (\gamma^2 - 1)}{q'},\end{aligned}\tag{13}$$

where  $y$  is the vector of past observations,  $var(\hat{X})/y$  is the conditional variance of  $\hat{X}$  given  $y$ , and  $E(\hat{X}/y)$  is the conditional expectation of  $\hat{X}$  given  $y$ . The variance of  $\hat{X}$  in equation (13) depends on  $X$  and increases with increasing levels of concentrations.

```
##### estimation of x #####
##### choose 1st observation of each lab as new observation #####

## 0 ug / L

## point estimate

y1 <- low[1,]
y2 <- low[6,]
y3 <- low[11,]

x <- function(y, alpha, beta, gamma) {
  xi_hat <- (y - alpha) / (beta * gamma)
}

x1 <- x(y1$Real, y1$alpha_i, y1$beta_i, gamma)
x2 <- x(y2$Real, y2$alpha_i, y2$beta_i, gamma)
x3 <- x(y3$Real, y3$alpha_i, y3$beta_i, gamma)
x_hat <- (x1 + x2 + x3) / 3
# x_hat
```

### 4.3 Variance and Confidence Interval of $X$

```
## variance

variance <- ((sigma_e2 / (y1$beta_i^2 * gamma^2)*(1 + 1 / 5)
  + sigma_e2 / (y2$beta_i^2 * gamma^2)*(1 + 1 / 5)
  + sigma_e2 / (y3$beta_i^2 * gamma^2)*(1 + 1 / 5)) / 3^2
  + 0^2 * (gamma^2 - 1) / 3)

## confidence interval
y_bar <- ( y1$Real + y2$Real + y3$Real ) / 3
sigma_alpha2 <- var(alpha_i)

c11 <- max(0, y_bar - 1.96*sqrt((sigma_e2 + sigma_alpha2) / 3))
clu <- y_bar + 1.96 * sqrt((sigma_e2 + sigma_alpha2) / 3)

CL_low <- round(c(c11, clu), 3)

# 20 ug / L

## point estimate
```

```

y12 <- twenty[1,]
y22 <- twenty[6,]
y32 <- twenty[11,]

x12 <- x(y12$Real, y12$alpha_i, y12$beta_i, gamma)
x22 <- x(y22$Real, y22$alpha_i, y22$beta_i, gamma)
x32 <- x(y32$Real, y32$alpha_i, y32$beta_i, gamma)
x_hat2 <- (x12 + x22 + x32) / 3
# x_hat2

## variance

variance2 <- ((sigma_e2 / (y12$beta_i^2 * gamma^2)*(1 + 1 / 5)
+ sigma_e2 / (y22$beta_i^2 * gamma^2)*(1 + 1 / 5)
+ sigma_e2 / (y32$beta_i^2 * gamma^2)*(1 + 1 / 5)) / 3^2
+ 20^2 * (gamma^2 - 1) / 3)

## confidence interval

X1 <- 20

c1i <- beta_i^2 * X1^2 * (gamma^4 - gamma^2) + sigma_e2
c2i <- c1i / (beta_i^2*X1^2)
c3i <- log((1 + sqrt(1 + 4*c2i)) / 2)

sum1 <- log(y12$Real - alpha_i[1]) / sqrt(c3i[1]) +
  log(y22$Real - alpha_i[2]) / sqrt(c3i[2]) +
  log(y32$Real - alpha_i[3]) / sqrt(c3i[3])

# lower bound

f1 <- function (x) {
  sum1 - log(beta_i[1] * x) / sqrt(c3i[1]) -
    log(beta_i[2] * x) / sqrt(c3i[2]) -
    log(beta_i[3] * x) / sqrt(c3i[3]) -
    1.96*sqrt(3)
}

# solve for x

c1l20 <- uniroot(f1, lower = 0.1, upper = 10000000)$root

# upper bound

```

```

f2 <- function (x) {
  sum1<- log(beta_i[1] * x) / sqrt(c3i[1]) -
    log(beta_i[2] * x) / sqrt(c3i[2]) -
    log(beta_i[3] * x) / sqrt(c3i[3]) +
    1.96*sqrt(3)
}

# solve for x

clu20 <- uniroot(f2, lower = 0.1, upper = 10000000)$root

CL_20 <- round(c(c1120, clu20), 3)

# 100 ug / L

## point estimate

y13 <- hundred[1,]
y23 <- hundred[6,]
y33 <- hundred[11,]

x13 <- x(y13$Real, y13$alpha_i, y13$beta_i, gamma)
x23 <- x(y23$Real, y23$alpha_i, y23$beta_i, gamma)
x33 <- x(y33$Real, y33$alpha_i, y33$beta_i, gamma)
x_hat3 <- (x13 + x23 + x33) / 3
# x_hat3

## variance

variance3 <- ((sigma_e2 / (y13$beta_i^2 * gamma^2)*(1 + 1 / 5)
  + sigma_e2 / (y23$beta_i^2 * gamma^2)*(1 + 1 / 5)
  + sigma_e2 / (y33$beta_i^2 * gamma^2)*(1 + 1 / 5)) / 3^2
  + 100^2 * (gamma^2 - 1) / 3)

## confidence interval

X2 <- 100

c1i2 <- beta_i^2 * X2^2 * (gamma^4 - gamma^2) + sigma_e2
c2i2 <- c1i2 / (beta_i^2*X2^2)
c3i2 <- log((1 + sqrt(1 + 4*c2i2)) / 2)

sum2 <- log(y13$Real - alpha_i[1]) / sqrt(c3i2[1]) +
  log(y23$Real - alpha_i[2]) / sqrt(c3i2[2]) +

```



```

log(y33$Real - alpha_i[3]) / sqrt(c3i2[3])

# lower bound

f3 <- function (x) {
  sum2- log(beta_i[1] * x) / sqrt(c3i2[1]) -
    log(beta_i[2] * x) / sqrt(c3i2[2]) -
    log(beta_i[3] * x) / sqrt(c3i2[3]) -
    1.96*sqrt(3)
}

# solve for x

c1l100 <- uniroot(f3, lower = 0.1, upper = 10000000)$root

# upper bound

f4 <- function (x) {
  sum2- log(beta_i[1] * x) / sqrt(c3i2[1]) -
    log(beta_i[2] * x) / sqrt(c3i2[2]) -
    log(beta_i[3] * x) / sqrt(c3i2[3]) +
    1.96*sqrt(3)
}

# solve for x

clu100 <- uniroot(f4, lower = 0.1, upper = 10000000)$root

CL_100 <- round(c(c1l100, clu100), 3)

x <- function(y, alpha, beta, gamma) {
  xi_hat <- (y - alpha) / (beta * gamma)
}

## simulated confidence level

scl_low <- numeric(1000)
scl_twenty <- numeric(1000)
scl_hundred <- numeric(1000)
alpha_i_sim <- matrix(NA, 1000, 5)
gamma_sim <- numeric(1000)
beta_i_sim <- matrix(NA, 1000, 5)
sigma_e2_sim <- numeric(1000)
x_hat_real <- numeric(1000)

```

```

x_hat_real_twenty <- numeric(1000)
x_hat_real_hundred <- numeric(1000)
for (i in 1 : 1000) {
  y_LOW<- matrix(NA, 5, 5)
  for (j in 1:5) { # lab
    for (k in 1:5) { # replication
      e_sim <- rnorm(1, 0, sqrt(sigma_e2))
      y_low <- alpha_i[j] + e_sim
      y_LOW[j, k] <- y_low
    }
  }
  y_twenty <- matrix(NA, 5, 5)
  for (j in 1:5) { # lab
    for (k in 1:5) { # replication
      e_sim2 <- rnorm(1, 0, sqrt(sigma_e2))
      eta_sim2 <- rnorm(1, 0, sqrt(sigma_eta2))
      y_twenty[j,k] <- alpha_i[j] + beta_i[j] * 20 * exp(eta_sim2) + e_sim2
    }
  }
  y_hundred <- matrix(NA, 5,5)
  for (j in 1:5) { # lab
    for (k in 1:5) { # replication
      e_sim3 <- rnorm(1, 0, sqrt(sigma_e2))
      eta_sim3 <- rnorm(1, 0, sqrt(sigma_eta2))
      y_hundred[j, k] <- alpha_i[j] + beta_i[j] * 100 * exp(eta_sim3) + e_sim3
    }
  }

  alpha_i_sim[i, ] <- apply(y_LOW, 1, mean) # alpha_i_sim
  sigma_e2_sim[i] <- sum(apply((y_LOW - alpha_i_sim[i, ])^2, 1, sum) / 4) / 5
  # sigma_e2_sim

  z_ijk_20 <- (y_twenty - alpha_i_sim[i,]) / 20
  z_ijk_100 <- (y_hundred - alpha_i_sim[i, ]) / 100

  mu_zij_20 <- apply(z_ijk_20, 1, sum) / 5
  mu_zij_100 <- apply(z_ijk_100, 1, sum) / 5
  mu_zi_sim <- (mu_zij_20 + mu_zij_100) / 2

  mu_zij_20_mat <- matrix(rep(mu_zij_20, 5), 5, 5, byrow = F)
  mu_zij_100_mat <- matrix(rep(mu_zij_100, 5), 5, 5, byrow = F)

  sigma_zi220_sim <- apply((z_ijk_20 - mu_zij_20_mat)^2 / 4, 1, sum)
  sigma_zi2100_sim <- apply((z_ijk_100 - mu_zij_100_mat)^2 / 4, 1, sum)

```

```

sigma_zi2_sim <- (sigma_zi220_sim + sigma_zi2100_sim) / 2

sigma_u2_sim <- sigma_e2_sim[i] * (1 / 20^2 + 1 / 100^2) / 2

beta_i_sim[i, ] <- sqrt(mu_zi_sim^4 / (sigma_zi2_sim - sigma_u2_sim + mu_zi_sim^2))

sigma_eta2_sim <- 2 * sum(log(mu_zi_sim / beta_i_sim[i,])) / 5

gamma_sim[i] <- exp(sigma_eta2_sim / 2)

labs <- sample(1:5, 3, F)
reps <- sample(1:5, 3, T)

m <- c(labs[1], reps[1])
n <- c(labs[2], reps[2])
o <- c(labs[3], reps[3])

y1_sim <- y_LOW[m[1], m[2]]
y2_sim <- y_LOW[n[1], n[2]]
y3_sim <- y_LOW[o[1], o[2]]

## simulated CI

y_bar_sim <- (y1_sim + y2_sim + y3_sim) / 3
sigma_alpha2_sim <- var(alpha_i_sim[i, ])

c11_sim <- max(0, y_bar_sim - 1.96 * sqrt((sigma_e2_sim[i] + sigma_alpha2_sim) / 3))
clu_sim <- y_bar_sim + 1.96 * sqrt((sigma_e2_sim[i] + sigma_alpha2_sim) / 3)

## Simulated X hat

x1_sim <- x(y1_sim, alpha_i_sim[i, m[1]], beta_i_sim[i, m[1]], gamma_sim)
x2_sim <- x(y2_sim, alpha_i_sim[i, n[1]], beta_i_sim[i, n[1]], gamma_sim)
x3_sim <- x(y3_sim, alpha_i_sim[i, o[1]], beta_i_sim[i, o[1]], gamma_sim)
x_hat_sim <- (x1_sim + x2_sim + x3_sim) / 3

## Real X hat
y1_real <- low[m[1]*m[2], ]$Real
y2_real <- low[n[1]*n[2], ]$Real
y3_real <- low[o[1]*o[2], ]$Real

x1_real <- x(y1_real, alpha_i[m[1]], beta_i[m[1]], gamma)
x2_real <- x(y2_real, alpha_i[n[1]], beta_i[n[1]], gamma)

```

```

x3_real <- x(y3_real, alpha_i[o[1]], beta_i[o[1]], gamma)
x_hat_real[i] <- (x1_real + x2_real + x3_real) / 3

if ((x_hat_real < clu_sim) & (x_hat_real > cll_sim)) {
  scl_low[i] <- 1
}

## Simulated Confidence interval for 20 ug / L concentration

y1_sim_twenty <- y_twenty[m[1], m[2]]
y2_sim_twenty <- y_twenty[n[1], n[2]]
y3_sim_twenty <- y_twenty[o[1], o[2]]

## Real X hat
y1_real_twenty <- twenty[m[1]*m[2], ]$Real
y2_real_twenty <- twenty[n[1]*n[2], ]$Real
y3_real_twenty <- twenty[o[1]*o[2], ]$Real

x1_real_twenty <- x(y1_real_twenty, alpha_i[m[1]], beta_i[m[1]], gamma)
x2_real_twenty <- x(y2_real_twenty, alpha_i[n[1]], beta_i[n[1]], gamma)
x3_real_twenty <- x(y3_real_twenty, alpha_i[o[1]], beta_i[o[1]], gamma)
x_hat_real_twenty[i] <- (x1_real_twenty + x2_real_twenty + x3_real_twenty) / 3

c1i_sim <- beta_i^2 * 20^2 * (gamma^4 - gamma^2) + sigma_e2
c2i_sim <- c1i_sim / (beta_i^2*20^2)
c3i_sim <- log((1 + sqrt(1 + 4*c2i_sim)) / 2)

sum1_sim <- log(y1_real_twenty - alpha_i[1]) / sqrt(c3i_sim[1]) +
  log(y2_real_twenty - alpha_i[2]) / sqrt(c3i_sim[2]) +
  log(y3_real_twenty - alpha_i[3]) / sqrt(c3i_sim[3])

# lower bound

f1_sim <- function (x) {
  sum1_sim - log(beta_i[1] * x) / sqrt(c3i_sim[1]) -
    log(beta_i[2] * x) / sqrt(c3i_sim[2]) -
    log(beta_i[3] * x) / sqrt(c3i_sim[3]) -
    1.96*sqrt(3)
}

# solve for x

c1l20_sim <- uniroot(f1_sim, lower = 0.1, upper = 10000000)$root

```

```

# upper bound

f2_sim <- function (x) {
  sum1_sim - log(beta_i[1] * x) / sqrt(c3i_sim[1]) -
    log(beta_i[2] * x) / sqrt(c3i_sim[2]) -
    log(beta_i[3] * x) / sqrt(c3i_sim[3]) +
    1.96*sqrt(3)
}

# solve for x

clu20_sim <- uniroot(f2_sim, lower = 0.1, upper = 10000000)$root

if ((x_hat_real_twenty < clu20_sim) & (x_hat_real_twenty > cll20_sim)) {
  scl_twenty[i] <- 1
}

## Simulated Confidence interval for 20 ug / L concentration

y1_sim_twenty <- y_twenty[m[1], m[2]]
y2_sim_twenty <- y_twenty[n[1], n[2]]
y3_sim_twenty <- y_twenty[o[1], o[2]]

## Real X hat
y1_real_twenty <- twenty[m[1]*m[2], ]$Real
y2_real_twenty <- twenty[n[1]*n[2], ]$Real
y3_real_twenty <- twenty[o[1]*o[2], ]$Real

x1_real_twenty <- x(y1_real_twenty, alpha_i[m[1]], beta_i[m[1]], gamma)
x2_real_twenty <- x(y2_real_twenty, alpha_i[n[1]], beta_i[n[1]], gamma)
x3_real_twenty <- x(y3_real_twenty, alpha_i[o[1]], beta_i[o[1]], gamma)
x_hat_real_twenty[i] <- (x1_real_twenty + x2_real_twenty + x3_real_twenty) / 3

cli_sim <- beta_i^2 * 20^2 * (gamma^4 - gamma^2) + sigma_e2
c2i_sim <- cli_sim / (beta_i^2*20^2)
c3i_sim <- log((1 + sqrt(1 + 4*c2i_sim)) / 2)

sum1_sim <- log(y1_real_twenty - alpha_i[1]) / sqrt(c3i_sim[1]) +
  log(y2_real_twenty - alpha_i[2]) / sqrt(c3i_sim[2]) +
  log(y3_real_twenty - alpha_i[3]) / sqrt(c3i_sim[3])

# lower bound

f1_sim <- function (x) {

```

```

sum1_sim - log(beta_i[1] * x) / sqrt(c3i_sim[1]) -
  log(beta_i[2] * x) / sqrt(c3i_sim[2]) -
  log(beta_i[3] * x) / sqrt(c3i_sim[3]) -
  1.96*sqrt(3)
}

# solve for x

c1l20_sim <- uniroot(f1_sim, lower = 0.1, upper = 10000000)$root

# upper bound

f2_sim <- function (x) {
  sum1_sim - log(beta_i[1] * x) / sqrt(c3i_sim[1]) -
    log(beta_i[2] * x) / sqrt(c3i_sim[2]) -
    log(beta_i[3] * x) / sqrt(c3i_sim[3]) +
    1.96*sqrt(3)
}

# solve for x

clu20_sim <- uniroot(f2_sim, lower = 0.1, upper = 10000000)$root

if ((x_hat_real_twenty < clu20_sim) & (x_hat_real_twenty > c1l20_sim)) {
  scl_twenty[i] <- 1
}

## Simulated Confidence interval for 100 ug / L concentration

y1_sim_hundred <- y_hundred[m[1], m[2]]
y2_sim_hundred <- y_hundred[n[1], n[2]]
y3_sim_hundred <- y_hundred[o[1], o[2]]

## Real X hat
y1_real_hundred <- hundred[m[1]*m[2], ]$Real
y2_real_hundred <- hundred[n[1]*n[2], ]$Real
y3_real_hundred <- hundred[o[1]*o[2], ]$Real

x1_real_hundred <- x(y1_real_hundred, alpha_i[m[1]], beta_i[m[1]], gamma)
x2_real_hundred <- x(y2_real_hundred, alpha_i[n[1]], beta_i[n[1]], gamma)
x3_real_hundred <- x(y3_real_hundred, alpha_i[o[1]], beta_i[o[1]], gamma)
x_hat_real_hundred[i] <- (x1_real_hundred + x2_real_hundred + x3_real_hundred) / 3

c1i2_sim <- beta_i^2 * 100^2 * (gamma^4 - gamma^2) + sigma_e2

```

```

c2i2_sim <- cli2_sim / (beta_i^2*100^2)
c3i2_sim <- log((1 + sqrt(1 + 4*c2i2_sim)) / 2)

sum2_sim <- log(y1_real_hundred - alpha_i[1]) / sqrt(c3i2_sim[1]) +
  log(y2_real_hundred - alpha_i[2]) / sqrt(c3i2_sim[2]) +
  log(y3_real_hundred - alpha_i[3]) / sqrt(c3i2_sim[3])

# lower bound

f1_sim_hundred <- function (x) {
  sum2_sim - log(beta_i[1] * x) / sqrt(c3i2_sim[1]) -
    log(beta_i[2] * x) / sqrt(c3i2_sim[2]) -
    log(beta_i[3] * x) / sqrt(c3i2_sim[3]) -
    1.96*sqrt(3)
}

# solve for x

c1l100_sim <- uniroot(f1_sim_hundred, lower = 0.1, upper = 10000000)$root

# upper bound

f2_sim_hundred <- function (x) {
  sum2_sim - log(beta_i[1] * x) / sqrt(c3i2_sim[1]) -
    log(beta_i[2] * x) / sqrt(c3i2_sim[2]) -
    log(beta_i[3] * x) / sqrt(c3i2_sim[3]) +
    1.96*sqrt(3)
}

# solve for x

clu100_sim <- uniroot(f2_sim_hundred, lower = 0.1, upper = 10000000)$root

if ((x_hat_real_hundred < clu100_sim) & (x_hat_real_hundred > c1l100_sim)) {
  scl_hundred[i] <- 1
}

}

scl_low <- sum(scl_low) / 1000
scl_20 <- sum(scl_twenty) / 1000
scl_100 <- sum(scl_hundred) / 1000

```

Variance for  $X$  is 3.905, 4.9507, 40.4201 for 0, 20, 100  $\mu\text{g/L}$  true concentration level,

respectively. 95% confidence interval for  $X$  is (0, 1.17), (15.493, 23.13), (90.767, 116.149), respectively.

Bhaumik and Gibbons [2005] proposed a method for constructing confidence regions, which is an approximation based on a normal or log-normal distribution.

Firstly, to construct a confidence region for low-level concentrations, let  $Y_{i0}$  be an observation collected from the  $i$ th laboratory with low-level true concentration. Define

$$\bar{Y}_0 = \sum_{i=1}^{q'} Y_{i0}/n_0.,$$

where  $n_0$  is the total number of measurements for low-level true concentration from all  $q'$  laboratories. For a low-level true concentration  $X_0$ , the  $(1 - \alpha)100\%$  confidence region of  $X_0$  is

$$(\max(0, \bar{Y}_0 - z_{\alpha/2}\sqrt{(\hat{\sigma}_e^2 + \hat{\sigma}_\alpha^2)/n_0}), \bar{y}_0 + z_{\alpha/2}\sqrt{(\hat{\sigma}_e^2 + \hat{\sigma}_\alpha^2)/n_0}).$$

$\hat{\sigma}_\alpha^2$  represents the variability of  $\alpha_i$  across all  $q$  laboratories in the calibration sample.

Secondly, to construct a confidence region for larger  $X$ , Bhaumik and Gibbons [2005] used the following lognormal approximation. Let

$$c_{1i} = \text{var}(Y_i) = \beta_i^2 X^2 (\gamma^4 - \gamma^2) + \sigma_e^2, c_{2i} = \frac{c_{1i}}{\beta_i^2 X^2}, \text{ and } c_{3i} = \ln\left(\frac{1 + \sqrt{1 + 4c_{2i}}}{2}\right).$$

Bhaumik and Gibbons [2005] also showed in a lemma that the approximate variance of  $\ln(\frac{Y_i - \alpha_i}{\beta_i X})$  is  $c_{3i}$ .

*Lemma 1.* Suppose that for an unknown concentration  $X$ , the corresponding observation  $Y_i$  collected from the  $i$ th laboratory follows model (2). Define  $V_i = \ln(\frac{Y_i - \alpha_i}{\beta_i X})$ . For a larger concentration  $X$ , the approximate variance of  $V_i$  is  $c_{3i}$ .

*Proof.* For a larger concentration  $X$ , the corresponding  $e_i$  in model (2) becomes insignificant compared to  $Xe^{\eta_i}$ , and hence the approximate distribution of  $(\frac{Y_i - \alpha_i}{\beta_i X})$  is lognormal. Using the expression for the variance of a lognormal distribution,  $c_{2i} = \text{var}(\frac{Y_i - \alpha_i}{\beta_i X})$  can be expressed as  $v^2 - v$ , for a positive number  $v$ . The positive root of the quadratic equation  $c_{2i} = v^2 - v$  is  $\frac{1 + \sqrt{1 + 4c_{2i}}}{2}$ , and hence the variance of  $\ln(\frac{Y_i - \alpha_i}{\beta_i X})$  is  $c_{3i}$ .

Let

$$Z_i(X) = \frac{\ln((Y_i - \alpha_i)/\beta_i X)}{\sqrt{c_{3i}}} = \frac{\ln(Y_i - \alpha_i) - \ln(\beta_i X)}{\sqrt{c_{3i}}}. \quad (14)$$

Thus  $Z_i(X) \sim N(0, 1)$  and the approximate distribution of  $Z(X) = \sum_{i=1}^{q'} Z_i(X)/\sqrt{q'} \sim N(0, 1)$ , where  $N(0, 1)$  denotes a standard normal distribution. We replace the parameters



on the right side of (14) by their corresponding estimates to compute their numerical values. Thus the  $(1 - \alpha)100\%$  confidence region for  $X$  is

$$\mathcal{R}(X) = \{X : -Z_{\alpha/2} \leq Z(X) \leq Z_{\alpha/2}\}. \quad (15)$$

## 5 Results and Discussions

The Original Results of the analysis of interlaboratory data for Cadmium given in Bhaumik and Gibbons [2005] are shown in the table below:

**Table: *Bhaumik and Gibbons [2005]'s Results***

True Concentration	X_hat	Var(X)	CI Lower Bound	CI Upper Bound	SCL
0 ug/L	0	2.76	0	1.182	0.945
20 ug/L	19.836	4.513	13.851	21.604	0.932
100 ug/L	97.377	45.241	85.197	111.506	0.954

This table displays point estimates, variances, and simulated confidence levels. According to their conclusions, this table reveals that

- The point estimates are close to the true concentrations
- Variances are proportional to concentrations
- All confidence intervals contain the true concentrations
- The width of the confidence intervals are proportional to concentration
- The simulated confidence levels are close to the intended confidence level of 95%.

Our results are summarized in a table using the same format as below:

**Table: *Our Results***

True Concentration	X_hat	Var(X)	CI Lower Bound	CI Upper Bound	SCL
0 ug/L	-1.577	3.905	0	1.17	0.938
20 ug/L	20.479	4.951	15.493	23.13	0.912
100 ug/L	102.137	40.42	90.767	116.149	0.982

We may notice that point estimate for  $\hat{X}$  at 0  $\mu\text{g}/L$  is negative, this is because we subjectively choose the first replicate from the first three laboratories, thus there may be some bias in estimation. But overall, we can see that our results are very close to theirs.

## References

- Dulal K Bhaumik and Robert D Gibbons. Confidence regions for random-effects calibration curves with heteroscedastic errors. *Technometrics*, 47(2):223–231, may 2005.
- David M. Rocke and Stefan Lorenzato. A two-component model for measurement error in analytical chemistry. *Technometrics*, 37(2):176–184, may 1995.