

Experimental Design and Analysis

HW06 Solution

Problem 1

(a)

The regression equations for location and dispersion, respectively, are

$$\begin{cases} \hat{y}_{\mathbf{x}} = 7.6360 + 0.1106x_B + 0.0519x_E + 0.0881x_C - 0.1298x_Q + 0.0423x_Bx_Q - 0.0827x_Cx_Q & (5.8) \\ \ln(\hat{\sigma}_{\mathbf{x}}^2) = -4.9313 + 0.9455x_B + 0.5556x_Dx_Q - 0.5445x_Bx_Cx_Q & (5.9) \end{cases}$$

We then choose factor settings $\mathbf{x} \in [-1, 1]^5$ to minimize the predicted mean squared error

$$\widehat{MSE}_{\mathbf{x}} = (\hat{y}_{\mathbf{x}} - 8)^2 + \exp(\ln(\hat{\sigma}_{\mathbf{x}}^2))$$

Using the *optim* function, we obtain

$$\arg \min_{\mathbf{x} \in [-1, 1]^5} \widehat{MSE}_{\mathbf{x}} = (-0.2723377, 1, 1, 1, -1)$$

(b)

Case	B	C	D	E	Q	$\hat{y}_{\mathbf{x}}$	$\hat{\sigma}_{\mathbf{x}}^2$	$\widehat{MSE}_{\mathbf{x}}$
Section 5.3	-1	1	1	2.5376	-1	8.000001	-6.976900	0.000933191
Exercise 1	-1	2.0749	1	-1	-1	7.999993	-7.562183	0.000519739
Exercise 2	-0.2723	1	1	1	-1	7.969902	-5.892627	0.003665612

To further reduce the predicted mean square error, we need to relax the range of the factor levels.

Factor Q is a good candidate to relax first because its main effect and its interaction terms have relatively large coefficients, and Q appears in both the location and dispersion models.

(c)

When we relax the range of the factor Q and allow x_Q to go below -1 (specifically, $x_Q \in [-1.6, 1]$ while keeping the other factors in $[-1, 1]$), the optimized setting improves further. Using the **optim** again, we obtain a new minimizer:

$$\mathbf{x} = (-1, 1, 1, -0.4080925, -1.6)$$

which yields a smaller predicted mean square error $\widehat{MSE}_{\mathbf{x}} \approx 0.0004822924$, outperforming the solutions obtained using the other methods.

However, this improvement comes from setting $x_Q = -1.6$, which lies outside the original design region $[-1, 1]$. Therefore, the result relies on extrapolation of the fitted regression models beyond the range of the experimental data, and the predicted gain should be interpreted cautiously (ideally confirmed by follow-up runs in the expanded region).

Problem 2

(a)

The defining words are $I = 134 = 235 = 1245$, so the resolution is III.

(b)

All the aliasing relations for this design are

I	= 134	= 235	= 1245
1	= 34	= 1235	= 245
2	= 1234	= 35	= 145
3	= 14	= 25	= 12345
4	= 13	= 2345	= 125
5	= 1345	= 23	= 124
12	= 234	= 135	= 45
15	= 345	= 123	= 24

We can see that the 2-factor interaction effects **13**, **14** and **34** are aliased with main effects. Therefore, only the interaction effects **12**, **23** and **24** can be estimated, provided that variable 5 is inert, all 2-factor interactions involving variable 5 and all higher order interactions are negligible.

Problem 3

(a)

We first write down the defining contrast subgroup for each choice of generators

$$\begin{cases} A: \mathbf{I} = \mathbf{12345} = \mathbf{1236} = \mathbf{456} \\ B: \mathbf{I} = \mathbf{1235} = \mathbf{2346} = \mathbf{1456} \end{cases}$$

From the defining relations, we find that design A is of resolution III, whereas design B is of resolution IV. Under the maximum resolution criterion, we recommend design B.

However, choosing the design with the highest resolution is not the only criterion. To gain a more detailed understanding of the alias structure, we next list the aliasing relations for designs A and B in turn, and then use these results to further justify our comparison and recommendation.

For design A, all the aliasing relations are

I	= 12345	= 1236	= 456
1^{SC}	= 2345	= 236	= 1456
2^{SC}	= 1345	= 136	= 2456
3^{SC}	= 1245	= 126	= 3456
4	= 1235	= 12346	= 56
5	= 1234	= 12356	= 46
6	= 123456	= 123	= 45
12	= 345	= 36	= 12456
13	= 245	= 26	= 13456
14^C	= 235	= 2346	= 156
15^C	= 234	= 2356	= 146
16	= 23456	= 23	= 145
24^C	= 135	= 1346	= 256
25^C	= 134	= 1356	= 246
34^C	= 125	= 1246	= 356
35^C	= 124	= 1256	= 346

For design B, all the aliasing relations are

I	= 1235	= 2346	= 1456
1^{SC}	= 235	= 12346	= 456
2^{SC}	= 135	= 346	= 12456
3^{SC}	= 125	= 246	= 13456
4^{SC}	= 12345	= 236	= 156
5^{SC}	= 123	= 23456	= 146
6^{SC}	= 12356	= 234	= 145
12	= 35	= 1346	= 2456
13	= 25	= 1246	= 3456
14	= 2345	= 1236	= 56
15	= 23	= 123456	= 46
16	= 2356	= 1234	= 45
24	= 1345	= 36	= 1256
26	= 1356	= 34	= 1245
124	= 345	= 136	= 256
126	= 356	= 134	= 245

After listing the aliasing relations, we summarize the clear and strongly clear effects for the two designs:

$$\begin{aligned}
 \text{Design A has } & \left\{ \begin{array}{ll} 3 \text{ clear main effects:} & \mathbf{1, 2, 3} \\ 6 \text{ clear 2-factor interactions:} & \mathbf{14, 15, 24, 25, 34, 35} \end{array} \right. \\
 \text{Design B has } & \left\{ \begin{array}{ll} 6 \text{ strongly clear main effects:} & \mathbf{1, 2, 3, 4, 5, 6} \\ 0 \text{ clear 2-factor interaction} & \end{array} \right.
 \end{aligned}$$

Therefore, if the objective is to obtain clear estimates of as many main effects as possible, design B is preferable (consistent with maximum resolution criterion). However, in applications where 2-factor interactions are of primary interest, design A can be more attractive, even though its overall resolution is lower.

(b)

A 2^{6-2} design has 16 runs, so the total degrees of freedom available for estimating factorial effects is 15. Now suppose, for contradiction, that there exists a 2^{6-2} design with the following two properties:

1. Main effects are clear
2. Two-factor interactions are clear

Consequently, each main effect and each two-factor interaction must belong to a distinct alias set. Hence, to estimate all main effects and all two-factor interactions separately, we would need at least

$$\binom{6}{1} + \binom{6}{2} = 21$$

degrees of freedom, which is impossible because the design provides only 15 degrees of freedom.

Therefore, a 2^{6-2} design does not exist.

Problem 4

(a)

All the aliasing relations of this experiment are

$$\begin{aligned}
 \mathbf{I} &= -ABCDE \\
 A^{SC} &= -BCDE \\
 B^{SC} &= -ACDE \\
 C^{SC} &= -ABDE \\
 D^{SC} &= -ABCE \\
 E^{SC} &= -ABCD \\
 AB^C &= -CDE \\
 AC^C &= -BDE \\
 AD^C &= -BCE \\
 AE^C &= -BCD \\
 BC^C &= -ADE \\
 BD^C &= -ACE \\
 BE^C &= -ACD \\
 CD^C &= -ABE \\
 CE^C &= -ABD \\
 DE^C &= -ABC
 \end{aligned}$$

From the aliasing relations above, we observe that all the main effects are strongly clear and all the 2-factor interactions are clear.

(b)

We list \bar{y} and $\ln(s^2)$ under each level combination below (If $s^2 = 0$, then we set $s = 0.001$).

```

data <- cbind(unique(data[,-6]), matrix(data[,6], ncol = 3, byrow = TRUE))
data[data == "-"] <- -1
data[data == "+"] <- 1
data <- as.data.frame(apply(data, 2, as.numeric))

y_bar <- apply(data[,6:8], 1, mean)
s_square <- apply(data[,6:8], 1, var); s_square[s_square == 0] <- 0.001
ln_s_square <- log(s_square)
data <- cbind(data[,1:5], y_bar, ln_s_square)

colnames(data) = c("A", "B", "C", "D", "E", "$\\bar{y}$", "ln$s^2$")
kable(data, digits = 3, row.names = F)

```

A	B	C	D	E	\bar{y}	$\ln s^2$
-1	-1	-1	-1	-1	1275.000	9.113
1	1	-1	-1	-1	1916.667	6.916
1	-1	1	-1	-1	1770.000	-6.908
-1	1	1	-1	-1	1275.000	-6.908
1	-1	-1	1	-1	1898.333	6.916
-1	1	-1	1	-1	1440.000	8.015
-1	-1	1	1	-1	1275.000	9.961
1	1	1	1	-1	2118.333	6.916
1	-1	-1	-1	1	1696.667	6.916
-1	1	-1	-1	1	1495.000	9.113
-1	-1	1	-1	1	1220.000	11.059
1	1	1	-1	1	1990.000	-6.908
-1	-1	-1	1	1	1495.000	10.580
1	1	-1	1	1	1641.667	6.916
1	-1	1	1	1	1916.667	6.916
-1	1	1	1	1	1256.667	6.916

In part (a), we know that all the main effects and all the 2-factor interactions are clear and hence estimable.

Thus, we fit the location model and the dispersion model, respectively, as follows:

$$\begin{cases} \bar{y} = \mathbf{X}\beta + \mathcal{E} \\ \ln(s^2) = \mathbf{X}\gamma + \mathcal{E}^* \end{cases}$$

where

$$\mathbf{X} = \begin{pmatrix} \text{(Int.)} & A & B & C & D & E & AB & AC & AD & AE & BC & BD & BE & CD & CE & DE \\ 1 & -1 & -1 & -1 & -1 & -1 & +1 & +1 & +1 & +1 & +1 & +1 & +1 & +1 & +1 & +1 \\ 1 & +1 & +1 & -1 & -1 & -1 & +1 & -1 & -1 & -1 & -1 & -1 & -1 & +1 & +1 & +1 \\ 1 & +1 & -1 & +1 & -1 & -1 & -1 & +1 & -1 & -1 & -1 & +1 & +1 & -1 & -1 & +1 \\ 1 & -1 & +1 & +1 & -1 & -1 & -1 & -1 & +1 & +1 & +1 & -1 & -1 & -1 & -1 & +1 \\ 1 & +1 & -1 & -1 & +1 & -1 & -1 & +1 & -1 & +1 & -1 & +1 & +1 & -1 & +1 & -1 \\ 1 & -1 & +1 & +1 & +1 & -1 & +1 & -1 & -1 & +1 & -1 & -1 & +1 & +1 & -1 & -1 \\ 1 & +1 & +1 & +1 & +1 & -1 & +1 & +1 & +1 & -1 & +1 & +1 & -1 & +1 & -1 & -1 \\ 1 & +1 & -1 & -1 & -1 & +1 & -1 & -1 & -1 & +1 & +1 & +1 & -1 & +1 & -1 & -1 \\ 1 & -1 & +1 & -1 & -1 & +1 & -1 & +1 & +1 & -1 & -1 & -1 & +1 & +1 & -1 & -1 \\ 1 & -1 & -1 & +1 & -1 & +1 & +1 & -1 & +1 & -1 & -1 & +1 & -1 & -1 & +1 & -1 \\ 1 & +1 & +1 & +1 & -1 & +1 & +1 & +1 & -1 & +1 & +1 & -1 & +1 & -1 & +1 & -1 \\ 1 & -1 & -1 & -1 & +1 & +1 & +1 & +1 & -1 & -1 & +1 & -1 & -1 & -1 & -1 & +1 \\ 1 & +1 & +1 & -1 & +1 & +1 & +1 & -1 & +1 & +1 & -1 & +1 & +1 & -1 & -1 & +1 \\ 1 & +1 & -1 & +1 & +1 & +1 & -1 & +1 & +1 & +1 & -1 & -1 & -1 & +1 & +1 & +1 \\ 1 & -1 & +1 & +1 & +1 & +1 & -1 & -1 & -1 & -1 & +1 & +1 & +1 & +1 & +1 & +1 \end{pmatrix}$$

consists of all the main effects and all the 2-factor interactions.

Then apply regression analysis to estimate the factorial effects:

```

y_bar <- data[,6]
ln_s_square <- data[,7]

location_model <- lm(
  y_bar ~ A + B + C + D + E +
    A:B + A:C + A:D + A:E +
    B:C + B:D + B:E +
    C:D + C:E +
    D:E,
  data = data
)

dispersion_model <- lm(
  ln_s_square ~ A + B + C + D + E +
    A:B + A:C + A:D + A:E +
    B:C + B:D + B:E +
    C:D + C:E +
    D:E,
  data = data
)

effect_name <- c("A", "B", "C", "D", "E", "AB", "AC", "AD", "AE", "BC", "BD", "BE", "CD", "CE", "DE")

location_effect <- 2 * coef(location_model)[-1]
dispersion_effect <- 2 * coef(dispersion_model)[-1]

names(location_effect) <- gsub(":", "", names(location_effect))
names(dispersion_effect) <- gsub(":", "", names(dispersion_effect))
location_effect <- location_effect[effect_name]
dispersion_effect <- dispersion_effect[effect_name]

table_effect <- cbind(
  effect_name,
  round(location_effect, 3),
  round(dispersion_effect, 3)
)

sig_loc <- c("A", "AC")
sig_dis <- c("C", "D", "CD")

fmt_cell <- function(val, is_sig, digits = 3){
  s <- sprintf(paste0("%.", digits, "f"), val)
  if (is_sig) paste0("$\\textcolor{red}{", s, "}")
  else paste0("$\\textcolor{black}{", s, "}")
}

loc_col <- mapply(fmt_cell, location_effect, effect_name %in% sig_loc,
  MoreArgs = list(digits = 3))
dis_col <- mapply(fmt_cell, dispersion_effect, effect_name %in% sig_dis,
  MoreArgs = list(digits = 3))

table_effect_colored <- cbind(effect_name, loc_col, dis_col)

```

```
colnames(table_effect_colored) <- c("Effect", "$\\bar y$", "$\\ln(s^2)$")
knitr::kable(table_effect_colored, row.names = FALSE, escape = FALSE)
```

Effect	\bar{y}	$\ln(s^2)$
A	527.083	-3.771
B	73.333	-2.947
C	-4.583	-5.430
D	50.417	5.093
E	-32.083	2.186
AB	22.917	2.947
AC	165.000	-1.482
AD	0.000	1.819
AE	-82.500	-2.186
BC	41.250	-2.306
BD	-105.417	1.545
BE	-59.583	-1.911
CD	27.500	5.001
CE	18.333	1.545
DE	-73.333	-2.306

Then check the significance of location and dispersion effects by *half-normal plots*:

```
stopifnot(exists("location_effect"), exists("dispersion_effect"))

if (is.null(names(location_effect))) {
  names(location_effect) <- paste0("E", seq_along(location_effect))
}
if (is.null(names(dispersion_effect))) {
  names(dispersion_effect) <- paste0("E", seq_along(dispersion_effect))
}

halfnorm_pink <- function(x, outlier = 1, main = "Half-normal Plot",
  hi_col = "deeppink",
  xlab = "Half-normal Quantiles",
  ylab = "Absolute Effects") {

  n <- length(x)

  ord <- order(abs(x))
  y <- abs(x)[ord]
  lab <- names(x)[ord]

  quan <- qnorm(0.5 + 0.5 * ((1:n) - 0.5) / n)

  xpad <- 0.08 * diff(range(quan))
  ypad <- 0.10 * diff(range(y))
  xlim <- c(min(quan) - xpad, max(quan) + xpad)
  ylim <- c(min(y) - ypad, max(y) + ypad)

  plot(quan, y, pch = 16, col = "black",
```

```

    main = main, xlab = xlab, ylab = ylab,
    xlim = xlim, ylim = ylim)

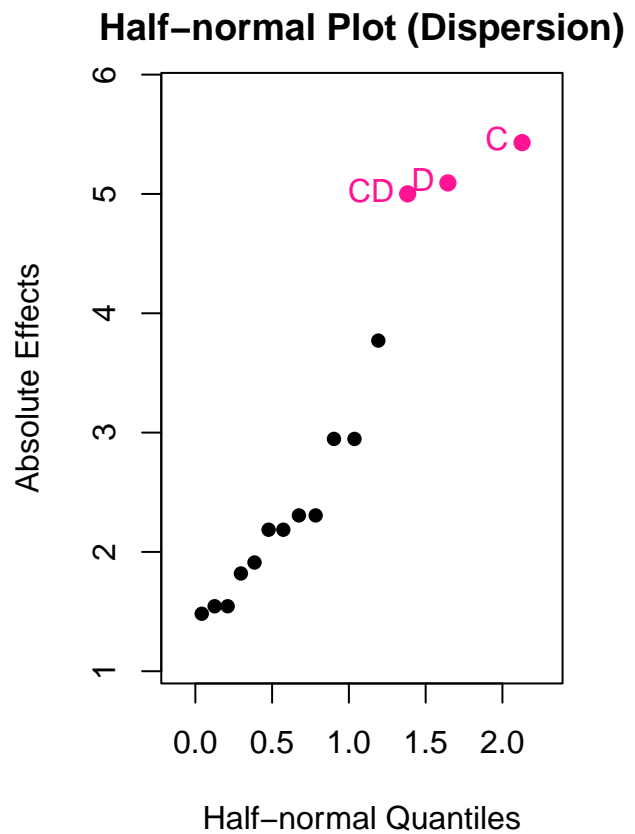
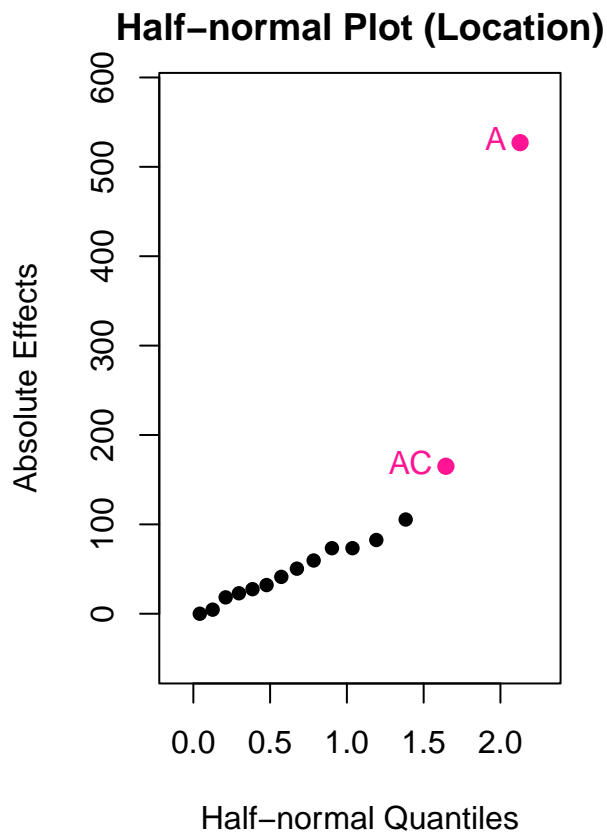
outlier <- min(outlier, n)
idx <- (n - outlier + 1):n

points(quan[idx], y[idx], pch = 16, col = hi_col, cex = 1.2)

pos_vec <- ifelse(quan[idx] > mean(xlim), 2, 4)
text(quan[idx], y[idx], labels = lab[idx],
     col = hi_col, pos = pos_vec, offset = 0.35, cex = 1, font = 1)
}

op <- par(mfrow = c(1, 2), mar = c(4.2, 4.2, 2.4, 1.6))
halfnorm_pink(location_effect, outlier = 2, main = "Half-normal Plot (Location)")
halfnorm_pink(dispersion_effect, outlier = 3, main = "Half-normal Plot (Dispersion)")

```



From the significance screening, effects *A* and *AC* are retained in the location model, while *C*, *D*, and *CD* are retained in the dispersion model. Therefore, the fitted models are:

$$\begin{cases} \hat{\bar{y}} = 1605 + 263.5415A + 82.5AC \\ \ln(\hat{s}^2) = 5.3456 - 2.715025C + 2.546284D + 2.500318CD \end{cases}$$

(c)

By inspecting the signs of the fitted dispersion model, we choose factor levels that minimize $\ln(\hat{s}^2)$. We set $C = +$ because its coefficient is negative, set $D = -$ because its coefficient is positive, and this choice also makes $CD = -$. This also makes $CD = -$ simultaneously, so the interaction term contributes negatively and further decreases $\ln(\hat{s}^2)$.

(d)

Since the coefficient of A is positive, we set $A = +$.

Since the coefficient of AC is also positive, we want $AC = +$, which (given $A = +$) implies $C = +$.

(e)

Because D is an adjustment factor and the response is larger-the-better, we adopt the following two-step procedure:

- Step 1: select $(A, C) = (+, +)$ so that

$$\hat{y} = 1605 + 263.5415 + 82.5 = 1951.$$

- Step 2: select $(C, D) = (+, -)$ so that

$$\hat{s}^2 = \exp\left(5.3456 - 2.715025 - 2.546284 - 2.500318\right) = 0.08927561.$$

That is, we first choose the factor levels to maximize the fitted mean, and then tune the adjustment factor to minimize the fitted variance.