

Incomplete Data Analysis Assignment (5%)

S1889112

A) (10 MARKS)

Carry out a CCA to find the mean value of the recovery time (and associated standard error) and to find also the correlations between the recovery time and the dose and between the recovery time and blood pressure.

A complete case analysis means dropping any missing data from the dataset, and then calculating our summary statistics.

Using a complete case analysis, the mean value of recovery time is 19.272, with an associated standard error of 2.442. The correlation between recovery time and the dose is 0.238, while the correlation between recovery time and blood pressure is -0.020. All results are to 3 decimal places.

B) (15 MARKS)

The same as in (a) but using mean imputation

Mean imputation takes the complete case mean, and uses this in place of the missing values.

Using mean imputation, the mean value of recovery time is unchanged at 19.272, but with an associated standard error of 2.284. The correlation between recovery time and the dose is 0.215, while the correlation between recovery time and blood pressure is -0.019. All results are to 3 decimal places.

C) (15 MARKS)

The same as in (a) but using mean regression imputation

Regression imputation takes the complete cases, performs linear regression (with our missing column as the response), and then uses the predictions $E(Y|X)$ in place of the missing data.

Using regression imputation, the mean value of recovery time is 19.444, with an associated standard error of 2.313. The correlation between recovery time and the dose is 0.279, while the correlation between recovery time and blood pressure is -0.011. All results are to 3 decimal places.

D) (15 MARKS)

The same as in (a) but using stochastic regression imputation. Do you need any extra care when conducting stochastic regression imputation in this example?

Stochastic regression imputation is identical to regression imputation but with an extra step. This involves adding noise to the predictions to try and restore lost variability from only using regression imputation. We assume the error is distributed normally with a mean of 0, and a homoscedastic variance equal to the estimated variance of the residuals from the complete case regression.

Using stochastic regression imputation, the mean value of recovery time is 19.903, with an associated standard error of 2.337. The correlation between recovery time and the dose is 0.295, while the correlation between recovery time and blood pressure is -0.0013. All results are to 3 decimal places.

Figure 1 shows the residual quantiles plotted against the quantiles for a normal curve (a qq plot). It shows that the sample quantiles for the residuals are consistently above the normal quantiles (illustrated by the lines). This suggests the residuals are not normally distributed, hence generating a normal error for the predictions is invalid.

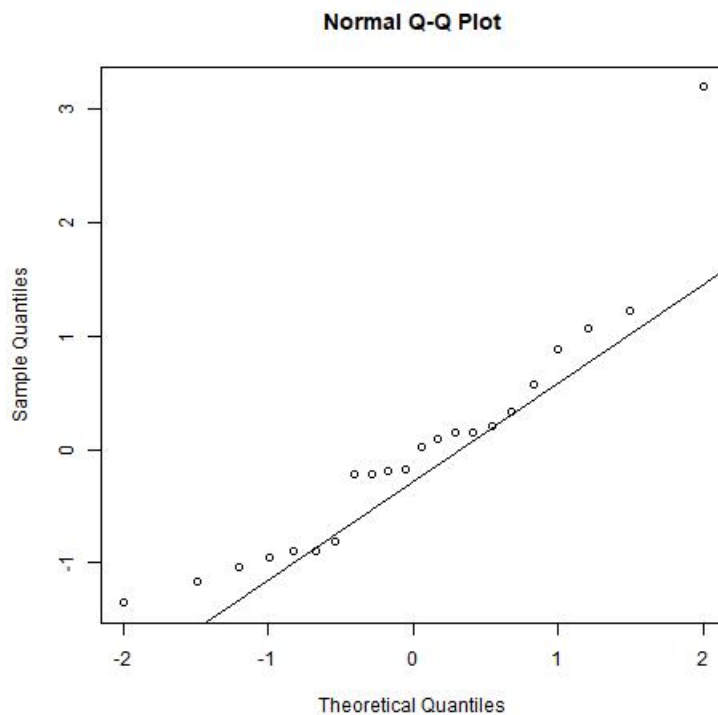


Figure 1: QQ plot

E) (30 MARKS)

You will now conduct the same analysis but applying another technique called predictive mean matching (Little, 1988), which is a special type of hot deck imputation. In the simplest form of this method (and the one you will use here), a regression model is used to predict the variables with missing values from the other (complete) variables. For each subject with a missing value, the donor is chosen to be the subject with a predicted value of her or his own that is closest (to be measured by the squared difference) to the prediction for the subject with the missing value.

Hot deck in general attempts to impute values by just using similar individuals instead of creating predictions. In the following example, by nature of linear regression - similar individuals implies similar predictions, and we use these predictions to find the most similar individual, and impute using their value.

Using hot deck imputation, the mean value of recovery time is 19.440, with an associated standard error of 2.464. The correlation between recovery time and the dose is 0.314, while the correlation between recovery time and blood pressure is -0.032. All results are to 3 decimal places.

F) (15 MARKS)

What is an obvious advantage of predictive mean matching over stochastic regression imputation?

Predictive mean matching avoids extrapolation beyond the range of values in the dataset, and so predictions based on this methodology are more likely to follow the true distribution. Additionally, it maintains discrete responses.

A R CODE

```
1 #-----#
2 ##### PREAMBLE #####
3 #-----#
4
5 load(file = 'databp.Rdata')
6 databp$missing <- ifelse(databp$R==0, 1, 0)
7 n <- nrow(databp)
8
9
10 #-----#
11 ##### PART A #####
12 #-----#
13
14 # Create df with only observed data
15
16 databp.cca <- databp[databp$missing == 0, ]
17
18
19 # Calculate statistics
20 mean.cca <- mean(databp.cca$recovtime)
21 se.cca <- sd (databp.cca$recovtime) / sqrt(n)
22 cor.recov.dose.cca <- cor (databp.cca$recovtime, exp(databp.cca$logdose))
23 cor.recov.bloodp.cca <- cor (databp.cca$recovtime, databp.cca$bloodp )
24
25
26 #### Results Part A ####
27
28 c(mean.cca, se.cca, cor.recov.dose.cca, cor.recov.bloodp.cca)
29
30
31
32 #-----#
33 ##### PART B #####
34 #-----#
35
36 # Already have complete case mean, which is the mean to impute
37 # Impute into new df
38 df.MI <- databp
39 df.MI$recovtime <- ifelse(df.MI$missing == 1, mean.cca, df.MI$recovtime)
40
41 # Calculate statistics
42
43
44 mean.MI <- mean(df.MI$recovtime)
45 se.MI <- sd (df.MI$recovtime) / sqrt(n)
46 cor.recov.dose.MI <- cor (df.MI$recovtime, exp(df.MI$logdose))
47 cor.recov.bloodp.MI <- cor (df.MI$recovtime, df.MI$bloodp)
48
49
50
51 #### Results Part B ####
52 c(mean.MI, se.MI, cor.recov.dose.MI, cor.recov.bloodp.MI)
53
54
55
```

```

56 #-----#
57 ##### PART C #####
58 #-----#
59
60 # Create new df
61 df.RI <- databp
62
63
64 # Fit model (lm ignores missing values), and generate predictions
65 mod <- lm(recovtime ~ bloodp + logdose, data = df.RI)
66 predictions.RI <- predict(mod, df.RI)
67
68 # Impute predictions to missing data
69 # If data missing, use prediction, otherwise use true value
70 df.RI$recovtime <- ifelse(df.RI$missing==1, predictions.RI, df.RI$recovtime)
71
72 # Calculate statistics
73
74
75 mean.RI <- mean(df.RI$recovtime)
76 se.RI <- sd (df.RI$recovtime)/sqrt(n)
77 cor.recov.dose.RI <- cor (df.RI$recovtime, exp(df.RI$logdose))
78 cor.recov.bloodp.RI <- cor (df.RI$recovtime, df.RI$bloodp)
79
80 ##### Results Part C #####
81 c(mean.RI, se.RI, cor.recov.dose.RI, cor.recov.bloodp.RI)
82
83
84 #-----#
85 ##### PART D #####
86 #-----#
87 set.seed(36)
88 # Already have regression results from part C
89 # Just need to add noise onto predictions
90
91 # SRI is just an extension of RI, so use df.RI as basis for new df
92 # Get noise scale parameter
93 df.SRI <- df.RI
94 noise.sd <- summary(mod)$sigma
95 # Add normal error if data was missing
96 df.SRI$recovtime <- ifelse(df.SRI$missing==1,
97                             df.SRI$recovtime + rnorm(1, 0, noise.sd),
98                             df.SRI$recovtime)
99
100 # Calculate statistics
101 mean.SRI <- mean(df.SRI$recovtime)
102 se.SRI <- sd (df.SRI$recovtime)/sqrt(n)
103 cor.recov.dose.SRI <- cor (df.SRI$recovtime, exp(df.SRI$logdose))
104 cor.recov.bloodp.SRI <- cor (df.SRI$recovtime, df.SRI$bloodp)
105
106 ##### Results Part D #####
107 c(mean.SRI, se.SRI, cor.recov.dose.SRI, cor.recov.bloodp.SRI)
108
109 # Linear regression assumption of normally distributed errors
110 jpeg('LaTeX/qplot.jpg')
111 qqnorm(rstandard(mod))
112 qqline(rstandard(mod))
113 dev.off()

```

```

113 #-----#
114 ##### PART E #####
115 #-----#
116 # Distance function
117 distance <- function(x,y) sqrt((x-y)^2)
118
119 # New df
120 df.HD <- databp
121
122 # Regression model recycled from RI section
123 mod.HD <- mod
124
125 # Predict for each item in data, and store in the dataframe for validity checks
126 df.HD$predictions <- predict.lm(mod.HD, newdata = df.HD)
127 df.HD$hotdeckpredictions <- numeric(n)
128
129 # Loop through rows
130 for (person in 1:n){
131   # If data not missing, just use observed value
132   if (df.HD[person, 'R'] == 1) {
133     df.HD[person, 'hotdeckpredictions'] <- df.HD[person, 'recovtime']
134   } else {
135     # Get missing person's predicted value, and calculate distance to
136     # all others' predicted value
137     predval <- df.HD[person, 'predictions']
138     distances <- distance(predval, df.HD$predictions)
139     # Find index of smallest distance that doesn't belong to the current or any
140     # other missing person
141     mindist <- min(distances[distances>0 & df.HD$missing==0])
142     minindex <- which(distances == mindist)
143     # Use that index's true value as the prediction for the missing person
144     df.HD[person, 'hotdeckpredictions'] <- df.HD[minindex, 'recovtime']
145     # Create min index column for manual inspection of process.
146     df.HD[person, 'minindex'] <- minindex
147   }
148 }
149
150 # View(df.HD)
151 # Sort by predictions in view to see method has used nearest prediction for
152 # the true value
153
154 # Perform imputation
155 df.HD$recovtime <- df.HD$hotdeckpredictions
156
157
158 # Calculate statistics
159
160 mean.HD <- mean(df.HD$recovtime)
161 se.HD <- sd (df.HD$recovtime)/sqrt(n)
162 cor.recov.dose.HD <- cor (df.HD$recovtime, exp(df.HD$logdose))
163 cor.recov.bloodp.HD <- cor (df.HD$recovtime, df.HD$bloodp)
164
165 #### Results Part E ####
166 c(mean.HD, se.HD, cor.recov.dose.HD, cor.recov.bloodp.HD)

```