# Calorie Consumption During Bicycle Work: A Statistical Analysis of an Incomplete Dataset

*Nuno Chicoria, Boris Shilov, Murat Cem Kose, Yibing Liu, Robin Vermote*

*19 March, 2018*

## Contents

## 1 Introduction

This project aimed to examine data originally gathered by Macdonald (1914) and conveyed to us by Greenwood and TF (1918), consisting of observations on seven people performing work using a bicycle ergometer, although our current dataset appears to include extra values and data not found in Greenwood and TF (1918), though these values may indeed be present in Macdonald (1914), access to which could not be obtained in a timely manner. In the body of this work it shall be assumed that every row in our dataset represents a separate individual, giving a total of 24 separate individuals across 24 rows. The dataset includes three separate measurements - weight of the individuals, calories per hour spent by individuals which serves as a measure of workout intensity, and calories spent during the task.

```
##    weight calhour calories
## 1    43.7    19.0       NA
## 2    43.7    43.0      279
## 3    43.7    56.0      346
## 4    54.6    13.0       NA
## 5    54.6    19.0       NA
## 6    54.6    43.0      280
## 7    54.6    56.0      335
## 8    55.7    13.0       NA
## 9    55.7    26.0      212
## 10   55.7    34.5      244
## 11   55.7    43.0      285
## 12   58.8    13.0       NA
## 13   58.8    43.0      298
## 14   60.5    19.0       NA
## 15   60.5    43.0      317
```

```
## 16   60.5    56.0      347
## 17   61.9    13.0       NA
## 18   61.9    19.0      216
## 19   61.9    34.5      265
## 20   61.9    43.0      306
## 21   61.9    56.0      348
## 22   66.7    13.0       NA
## 23   66.7    43.0      324
## 24   66.7    56.0      352
```

# 2   Methods and Results

## 2.1   Data exploration

A set of summary statistics for the dataset is presented below. It can be immediately seen that the response calories variable is missing in eight cases and is the only incomplete variable in the dataset. The mean and median values for all variables in the dataset are very similar to each other which indicates a symmetric distribution. A matrix of summary plots for the dataset is presented in Figure 2. We can clearly see what appears to be an extremely strong positive correlation between calories and workout intensity (0.95), and a very small positive correlation between calories and weight(0.11). The scatterplot further indicates that the correlation between calories and workout intensity is very likely to be linear.

The distributions of the values are plotted as boxplots in Fig. 1. Note that the response variable is plotted with missing values excluded in all of these figures, thus despite expecting an approximately similar distribution between workout intensity and calories variables, the calories distribution is shifted upwards due to the missing values.

```
##                  weight  calhour  calories
## nbr.val          24.0000  24.0000   16.0000
## nbr.null          0.0000   0.0000    0.0000
## nbr.na            0.0000   0.0000    8.0000
## min              43.7000  13.0000  212.0000
## max              66.7000  56.0000  352.0000
## range            23.0000  43.0000  140.0000
## sum            1381.0000 817.0000 4754.0000
## median           58.8000  38.7500  302.0000
## mean             57.5417  34.0417  297.1250
## SE.mean           1.3453   3.3396   11.4669
## CI.mean.0.95      2.7829   6.9085   24.4412
## var              43.4338 267.6721 2103.8500
## std.dev           6.5904  16.3607   45.8677
## coef.var          0.1145   0.4806    0.1544
```

We check the signifance of the two positive correlations we have found using Pearson's correlation (using Central Limit Theorem as the dataset contains around 20 rows). Here $H_0 : correlation = 0$; $H1 : correlation \neq 0$; $95\%CI$.

```
##
##  Pearson's product-moment correlation
##
## data:  muscledata_edit$calhour and muscledata_edit$calories
## t = 12, df = 14, p-value = 2e-08
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
```
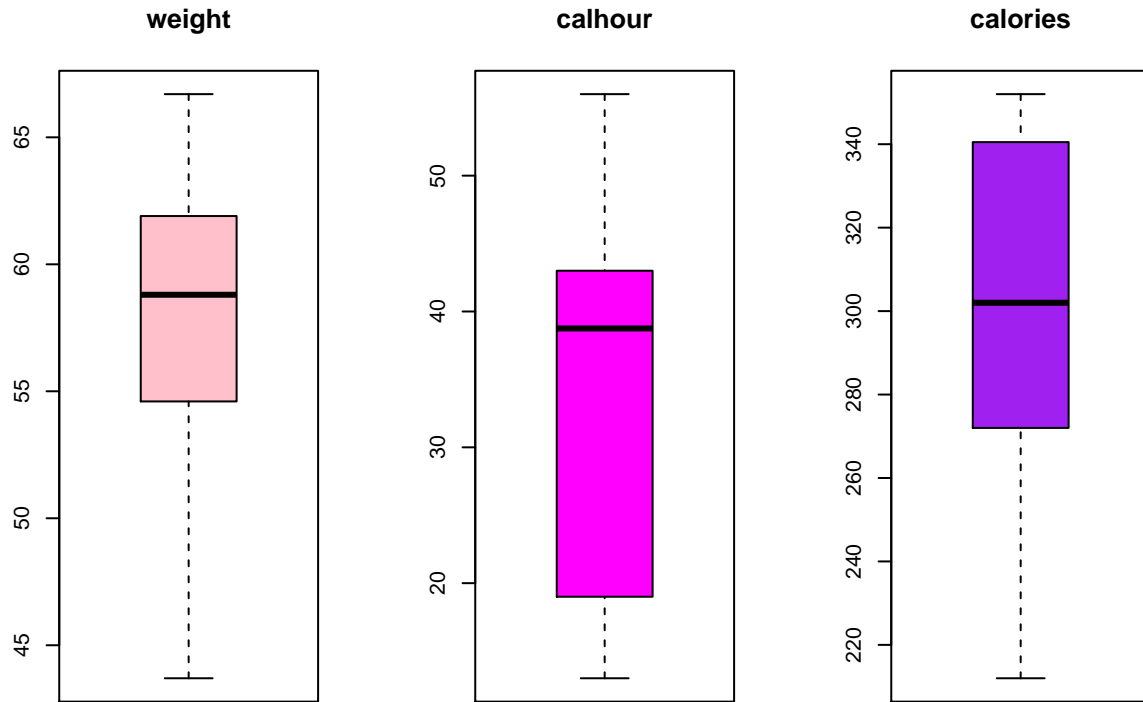
Figure 1: Boxplots for the dependent variables weight, calhour and independent variable calories.
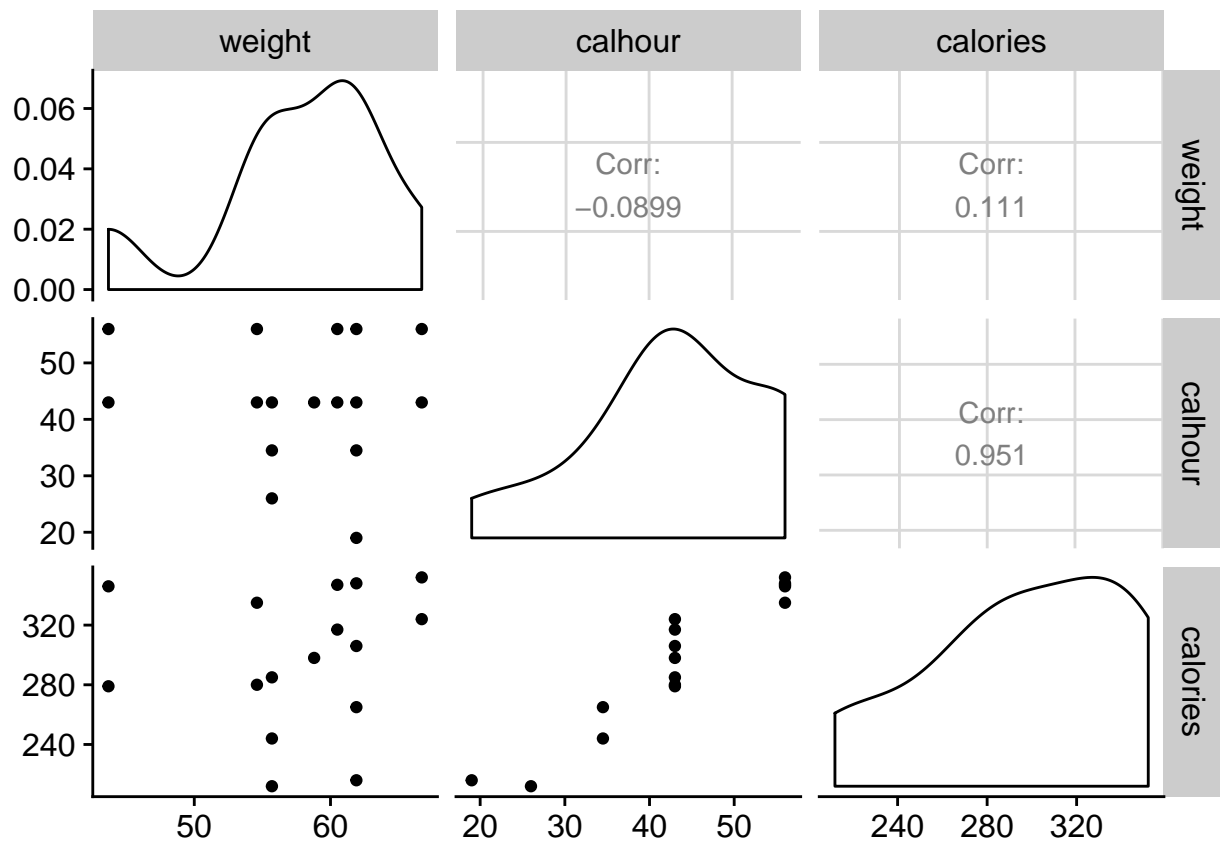


Figure 2: A summary statistics plot of the dataset using the ggplot command.
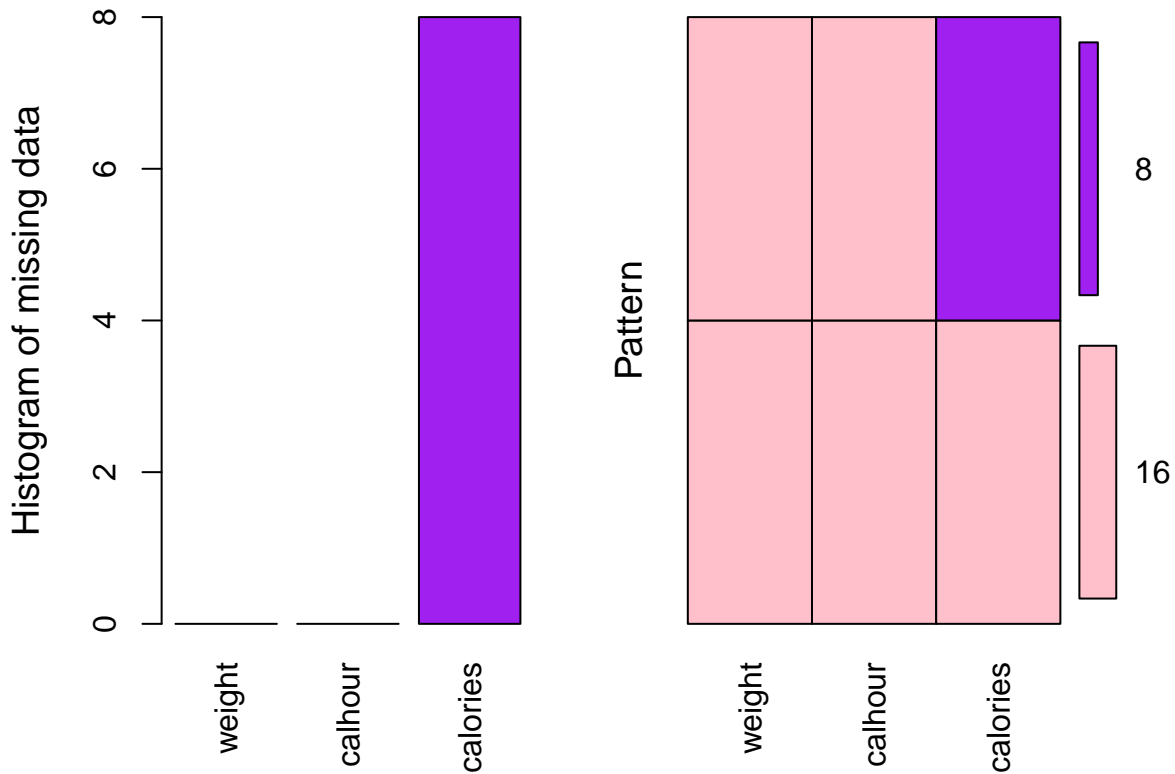
Figure 3: Patterns of missing data across variables.

```
##  0.8615 0.9832
## sample estimates:
##    cor
## 0.9511

##
##  Pearson's product-moment correlation
##
## data:  muscledata_edit$weight and muscledata_edit$calories
## t = 0.42, df = 14, p-value = 0.7
## alternative hypothesis: true correlation is not equal to 0
## 95 percent confidence interval:
##  -0.4068  0.5753
## sample estimates:
##    cor
## 0.1114
```

Clearly we reject the null hypothesis with regards to workout intensity and calories and accept it with regards to weight and calories. This indicates that there is a non-spurious correlation between workout intensity and calories in the population.

## 2.2  Missing data exploration

A histogram of missing data is shown in Fig. 3. We confirm our previous observation that all the missing values are located in our response variable.

Fig. 4A and C we see that the missing data approximately evenly distributed among the different weight

4

variables. In Fig. 4B and D we see that the missing data distribution is extremely biased towards the lower end of the range with regards to workout intensity. This may be because of the difficulty in measuring heat production at lower exercise intensity - in other words, the missingness is likely systematic due to technical noise. Importantly, the missingness appears to depend only on an observed variable in this study - the calories. Thus, this suggests "Missing-at-Random" as the most probable missing data mechanism, allowing us to proceed with applying missing data strategies - particularly MI and IPW. We will nonetheless evaluate some of the more common methods as well.

## 2.3 Complete case analysis

Complete case analysis relies on removing rows of our dataset that have missing values - giving us a restricted sample of 16 rows to work with.

We select the best linear model to use for complete case analysis using stepwise Akaike's Information Criterion - a measure of the quality of our statistical models relative to each other, which indicates the amount of information lost by excluding or including model terms.

```
## Start:  AIC=123.4
## calories ~ 1
##
##            Df Sum of Sq   RSS    AIC
## + calhour  1     28544  3014   87.8
## <none>                 31558  123.4
## + weight   1       392 31166  125.2
##
## Step:  AIC=87.81
## calories ~ calhour
##
##            Df Sum of Sq   RSS    AIC
## + weight   1      1234  1780   81.4
## <none>                  3014   87.8
## - calhour  1     28544 31558  123.4
##
## Step:  AIC=81.39
## calories ~ calhour + weight
##
##                   Df Sum of Sq   RSS    AIC
## + weight:calhour  1       782   998   74.1
## <none>                         1780   81.4
## - weight          1      1234  3014   87.8
## - calhour         1     29386 31166  125.2
##
## Step:  AIC=74.13
## calories ~ calhour + weight + calhour:weight
##
##                   Df Sum of Sq  RSS   AIC
## <none>                          998  74.1
## - calhour:weight  1       782 1780  81.4
```

Lower AIC is better, thus we conclude that a model incorporating weight, calhour and an interaction term is the most explanatory linear model available given our exploratory analysis. In mathematical terms:

$$calories_i = \beta_0 + \beta_1 * weight_i + \beta_2 * calhour_i + \beta_3 * (weight_i * calhour_i) + \epsilon_i$$
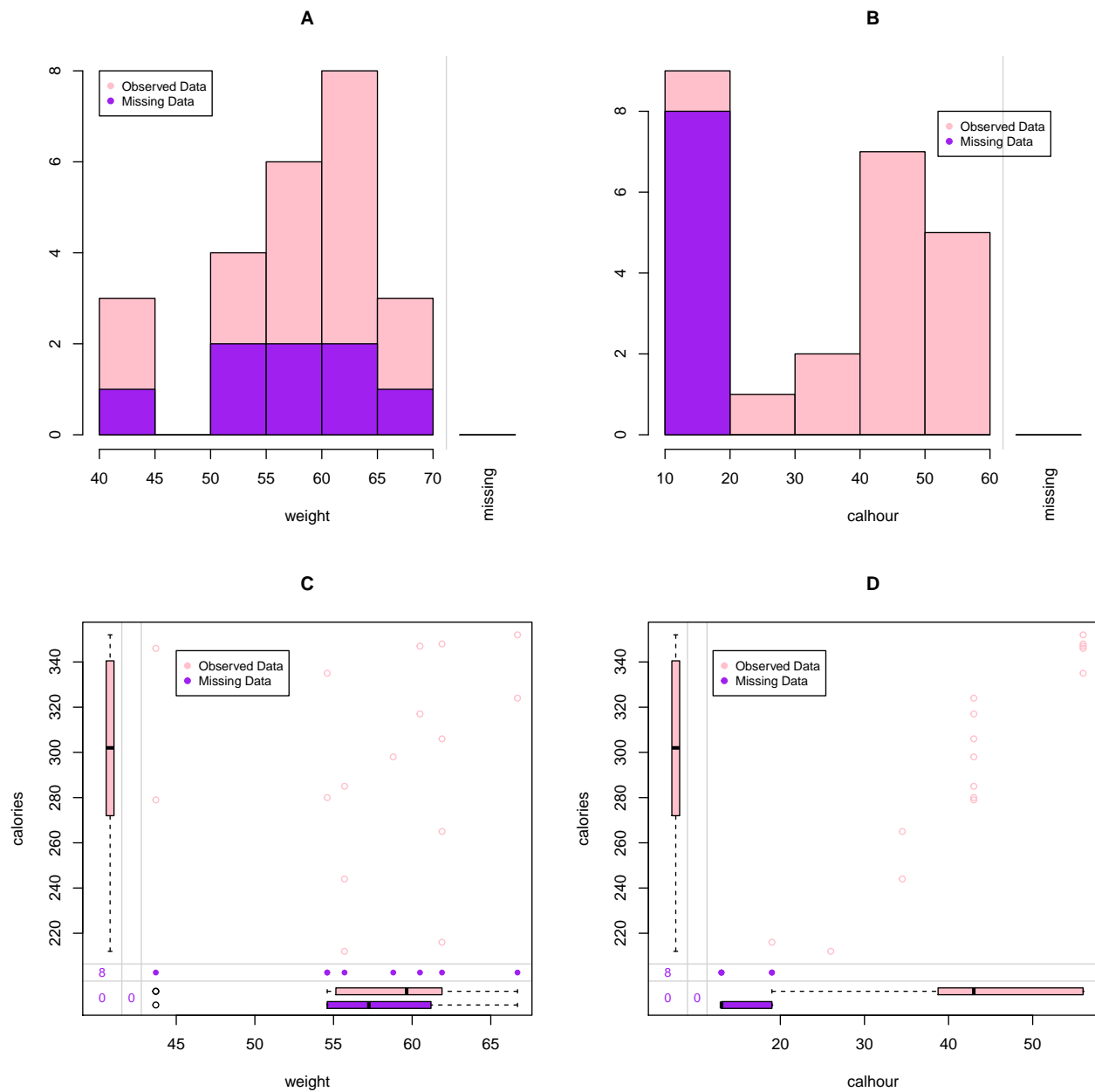
The summary for this model:

5

Figure 4: Histograms of the observed and missing data as well as marginplots depicting histograms and correlations.

```
## 
## Call:
## lm(formula = calories ~ weight + calhour + weight * calhour,
##     data = muscledata_edit)
## 
## Residuals:
##    Min     1Q Median     3Q    Max
## -12.48  -5.70  -1.04   2.39  16.95
## 
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -330.884    124.674   -2.65  0.02102 *
## weight           7.728      2.106    3.67  0.00321 **
## calhour         11.787      2.548    4.63  0.00058 ***
## weight:calhour  -0.132      0.043   -3.07  0.00977 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
## 
## Residual standard error: 9.12 on 12 degrees of freedom
## Multiple R-squared:  0.968,  Adjusted R-squared:  0.96
## F-statistic:  123 on 3 and 12 DF,  p-value: 2.89e-09
```

All four terms appear highly significant in this model. However, the intercept term $\beta_0$ does not have a physical meaning in this model. The significance of the interaction term means that the weight has an influence on the effect of workout intensity on calories.

An effects plot is presented in Fig. 6. This plot indicates that there is a decrease in slope, determined by workout intensity and calories, as workout intensity is increasing.

## 2.4   Multiple imputation analysis

Multiple imputation is an approach to deal with incomplete data that can be applied to univariate or multivariate data. The technique replaces missing values with two or more imputed values. Unlike simpler single imputation methods where only a single value is imputed, such as mean imputation, multiple imputation as the name suggests replaces each missing value with multiple imputed values, in effect generating a number of datasets. Practically, this allows us to represent a variety of theoretical mechanisms for why the nonresponse occured. These differing datasets are known as multiply imputed datasets. These datasets are used to generate a matrix of regression coefficients, in essence building a regression model. We generate as many regression coefficient matrices as there are multiply imputed datasets. We then pool the regression coefficients into a single estimate which can be used to estimate variance (Rubin 2004). There are several methods of imputation available.

First we use the predictive mean matching (PMM) method. This is one of the "default" methods and it faithfully reproduces the relations present in the original data even if they happen to be nonlinear. The results of such a simulation with 100 imputed value datasets:

According to the resulting PMM model, none of the possible dependent variables have any statistically significant influence on our response variable. The effects plot in Fig 6 shows that there is no influence of the interaction term on the response variable since the slope does not change. Notice also the contraction of the 95% confidence limit due to the imputation process.

The strip plot is shown in Fig. 7 showing original data in pink and generated data in purple. We thus indeed confirm PMM-generated data follows very similar relations to the original dataset.

An alternative to PMM with our dataset is Bayesian linear regression, which can create univariate missing data. This reformulates our linear model in probabilistic terms. In this method we assume a prior distribution
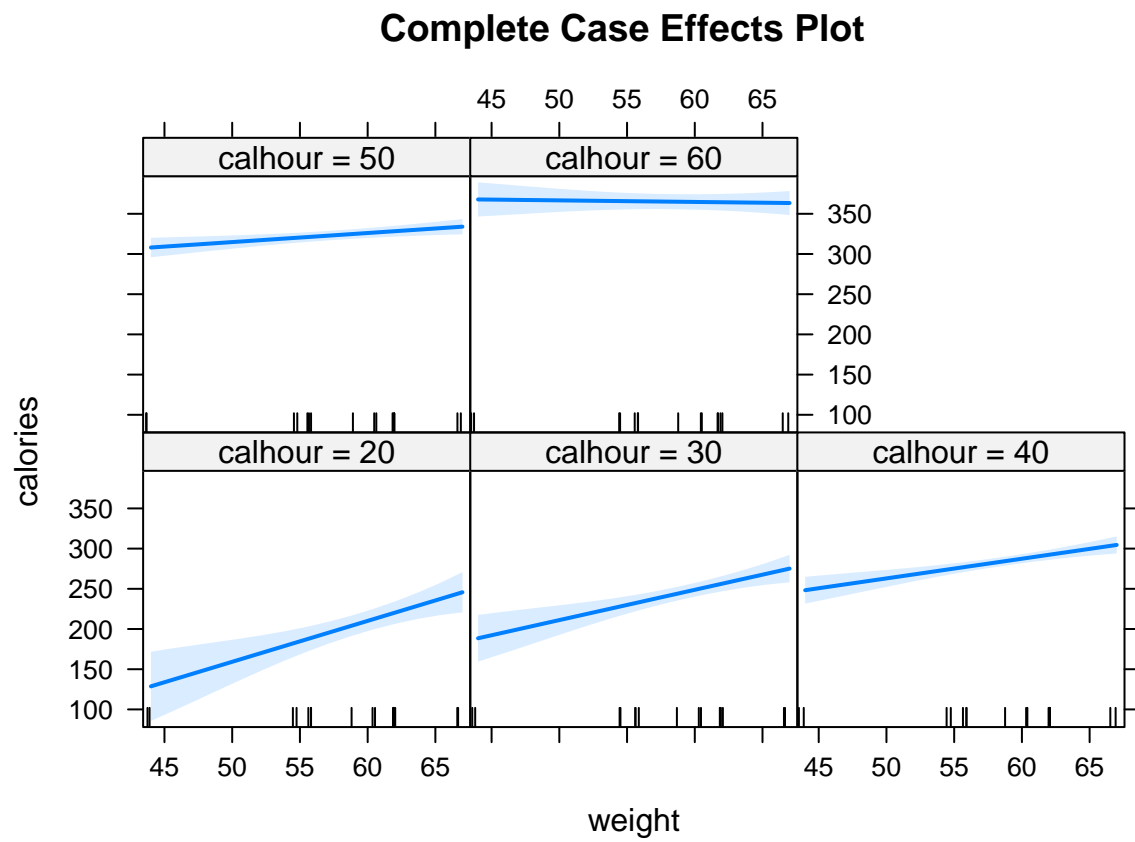
**Complete Case Effects Plot**


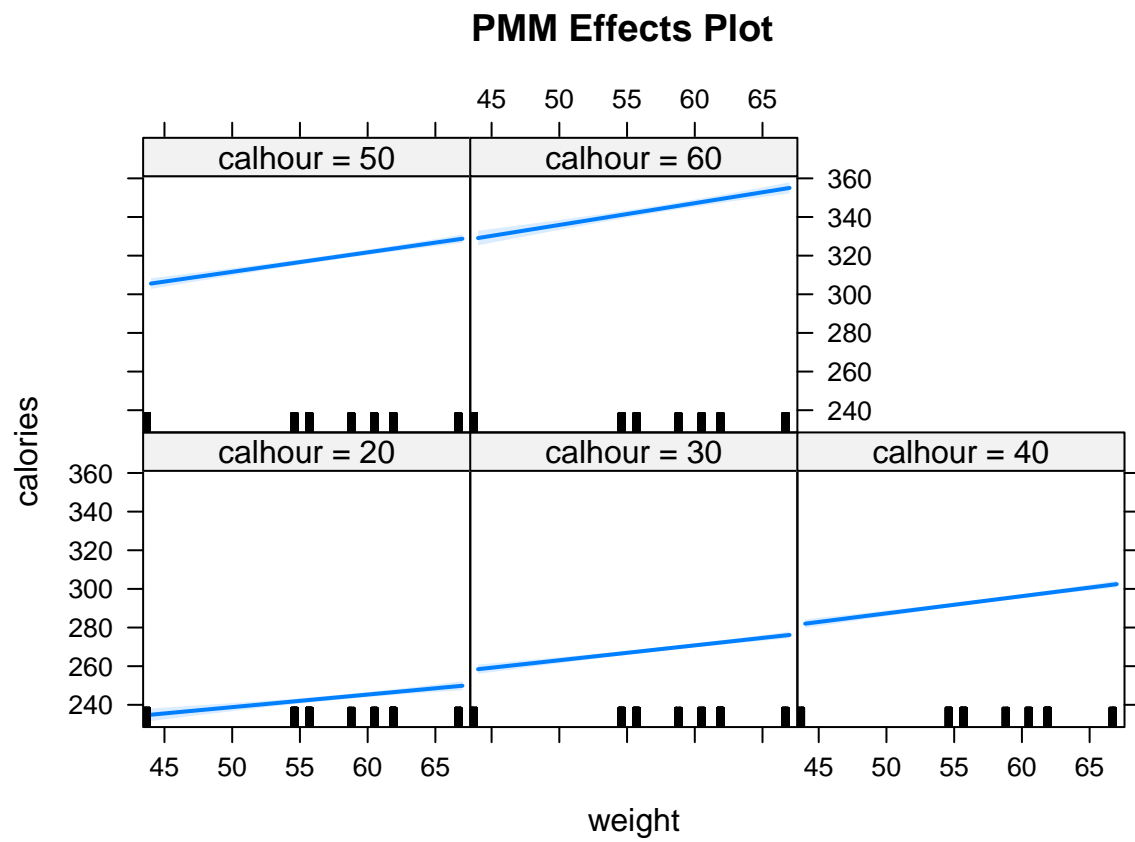
Figure 5: The All Effects plot for the Complete Case linear model.

Figure 6: The All Effects plot for MI using the PMM method.

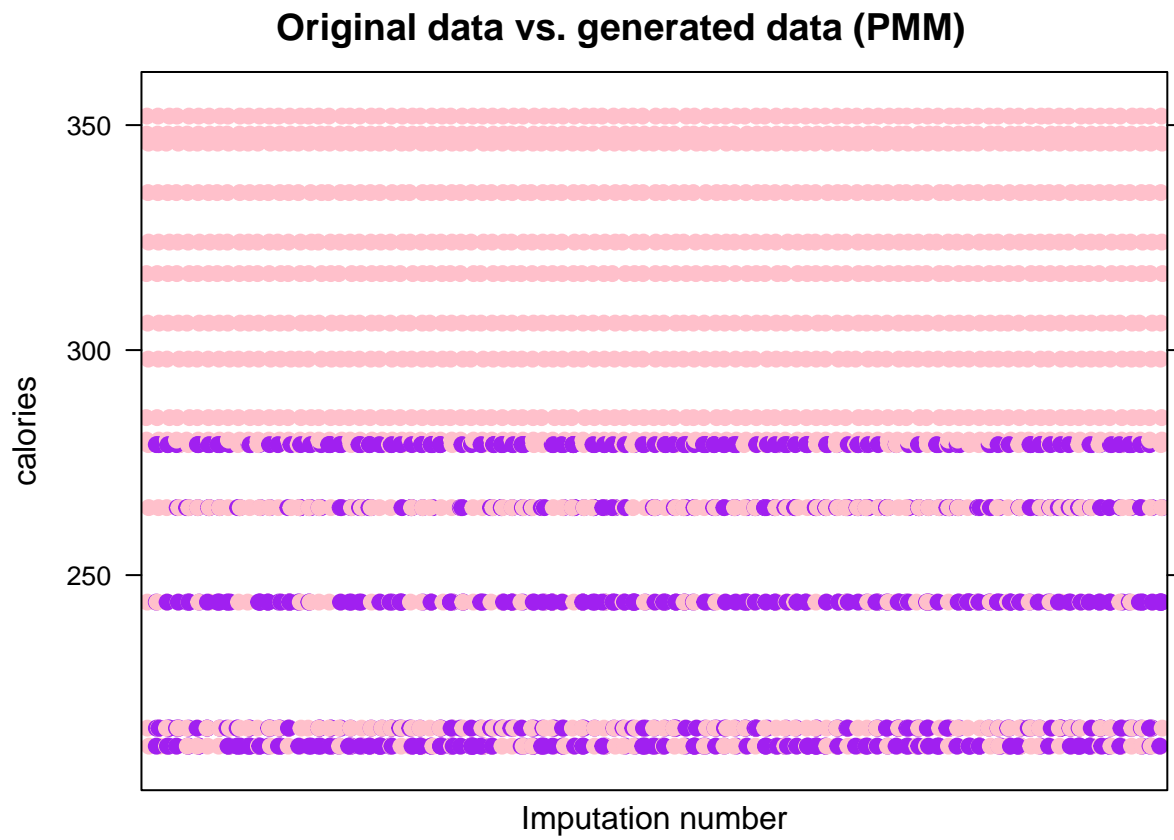Figure 7: The strip plot of PMM data.
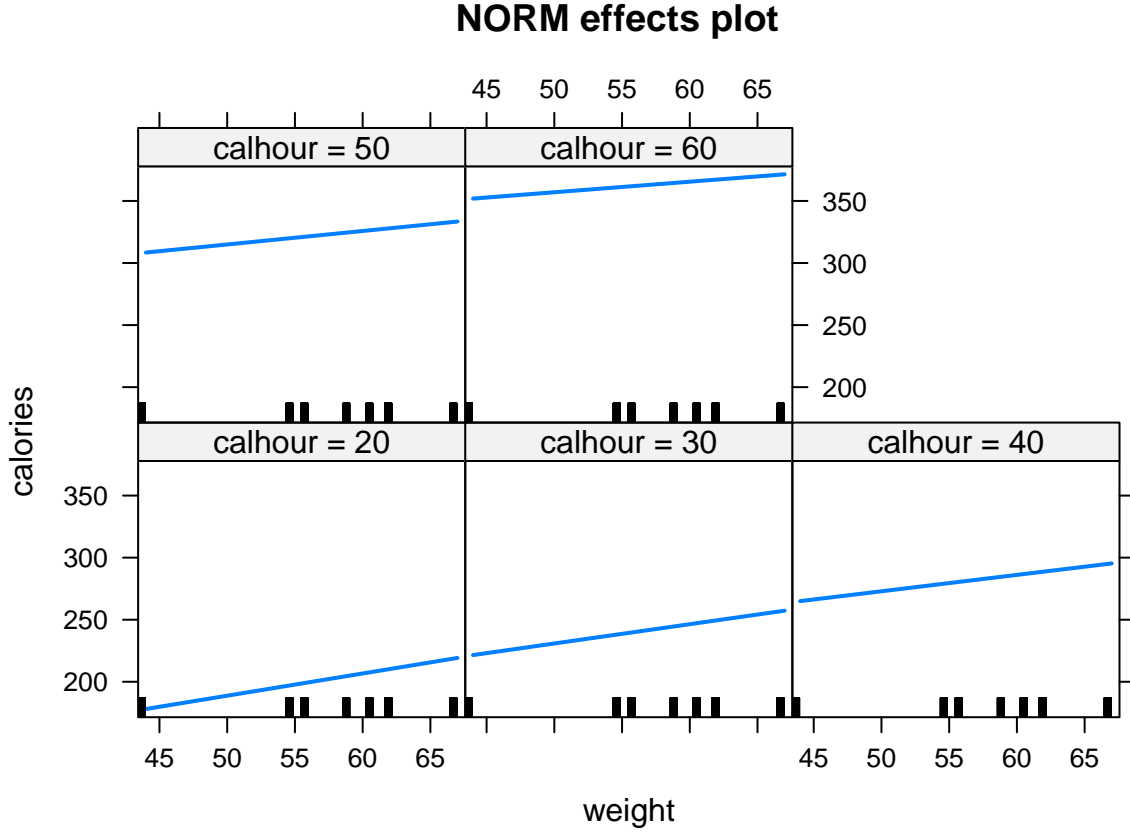
**NORM effects plot**



Figure 8: The All Effects plot for MI using the Bayesian NORM method.

for the parameters of the regression (Rubin 2004). In our method the prior distribution is assumed to be normal (the normal model). The result is:

```
##                     est       se        t    df Pr(>|t|)     lo 95
## (Intercept)    -7.93378 90.91732 -0.08726 6.318  0.93316 -227.7140
## weight          2.25300  1.51549  1.48665 6.706  0.18254    -1.3627
## calhour         5.37306  1.95778  2.74447 8.052  0.02512     0.8634
## weight:calhour -0.02337  0.03273 -0.71390 8.491  0.49444    -0.0981
##                   hi 95 nmis    fmi lambda
## (Intercept)    211.84642   NA 0.7209 0.6446
## weight           5.86870    0 0.7006 0.6229
## calhour          9.88268    0 0.6298 0.5480
## weight:calhour   0.05136   NA 0.6067 0.5238
```

According to our Bayesian model, the workout intensity is the only variable that has a statistically significant influence on our response variable. The effects plot in Fig. 8 demonstrates a similar finding to the PMM effects plot in that due to preserving the slope remaining the same the interaction variable clearly does not have much effect. The strip plot in Fig. 9 highlights the probabilistic nature of the Bayesian imputation algorithm, the imputed data points being sampled from the obtained posterior distribution of our parameters.

## 2.5 Inverse Probability Weighting analysis

The Inverse Probability Weighting method attemps to mitigate the bias introduced by complete case studies if the excluded population appears to be systematically different from the complete cases. The complete cases
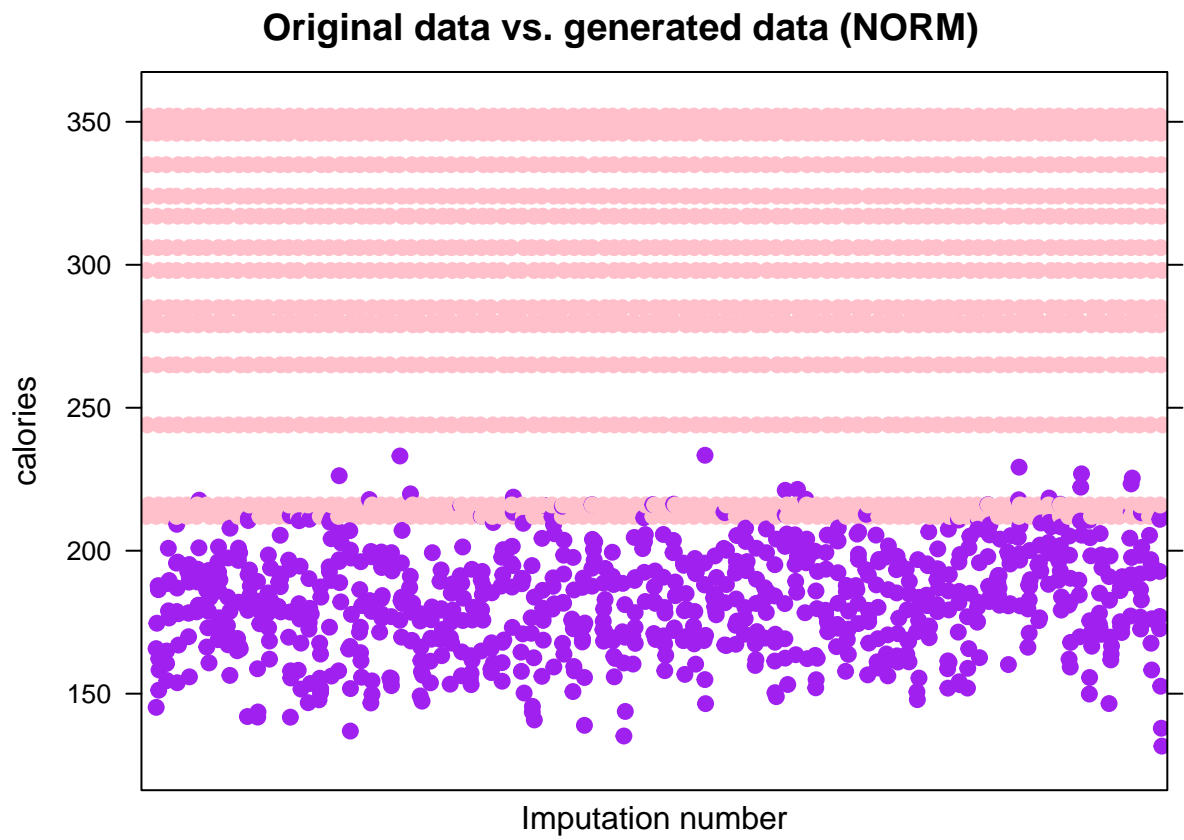
Figure 9: The strip plot of Bayesian NORM data.

are weighted using the inverse probability of their being a complete case. As you may recall, this is intuitively a very plausible model for our data since a mechanistic hypothesis for the MAR present is due to technical noise. To emphasize, in IPW the analytical model is only fitted to complete cases. The results are thus:

```
## Warning: In lm.fit(x, y, offset = offset, singular.ok = singular.ok, ...) :
##  extra argument 'family' will be disregarded

##
## Call:
## lm(formula = calories ~ weight + calhour + weight * calhour,
##     data = IPWanal_muscledata, weights = muscledata$w)
##
## Weighted Residuals:
##    Min     1Q Median     3Q    Max
##  -91.0  -40.5  -11.0   20.1  129.8
##
## Coefficients:
##               Estimate Std. Error t value Pr(>|t|)
## (Intercept)   -353.7928   129.1577   -2.74  0.01796 *
## weight           8.1131     2.1698    3.74  0.00283 **
## calhour         12.1321     2.6513    4.58  0.00064 ***
## weight:calhour  -0.1378     0.0445   -3.10  0.00926 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 68.2 on 12 degrees of freedom
##   (8 observations deleted due to missingness)
## Multiple R-squared:  0.97,   Adjusted R-squared:  0.962
## F-statistic:  128 on 3 and 12 DF,  p-value: 2.25e-09
```

All the parameters are highly significant, similarly to the complete cases analysis as can be expected. Fig. 10 effects plot is also highly similar.

The relative AIC values of the complete case and IPW models can be used for comparison to validate that our IPW model is indeed better than simple CC:

```
## [1] 121.5
```

```
## [1] 121.1
```

We can observed that IPW yields only a miniscule improvement over CC in this case.

# 3   Discussion

Due to the NA values, we conducted a full model analysis with a complete case and three NA comparsions (you can write this better) models. Beacuse the NA values are not evenly distrubited among calhour, we decided to try different approaches for NA handling.

IPW assigns weights to each observation so it uses already availible ones. Since all calories values in calhour 13 are missing, the method cannot assign a weight. no value can represent this group, other missing values fall into calhour 19, while a higher weight is assigned to the only available data in calhour 19. so the only difference between cc and ipw is only based on this value, thus the graph is the mostly the same for both CC and IPW and that's why we chose to represent both with the same graph.

for the MI method, we tried two ways to generate the data corresponding to the missing data, PMM and NORM. PMM generates the data according to the pattern in the observed ones. in our cases, the data is
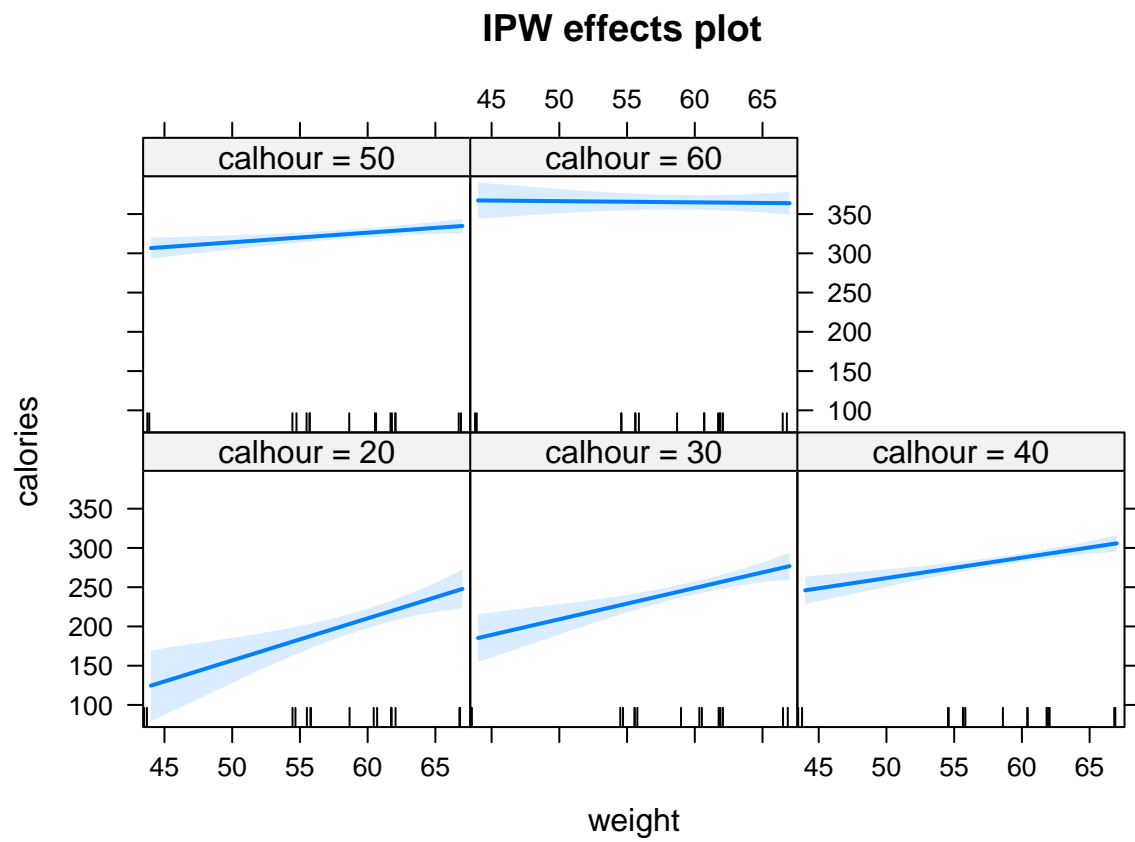
13

Figure 10: The All Effects plot for our IPW-modelled data.

discreted by the body weıght, so pmm generated the data dıscreted as well. ın norm method, the data ıs generated based on normal dıstrıbutıon.

In our MI models, there ıs no changes ın slope between each calhour and weıght ınteractıons. by checkıng the fıttıng model generated, the p-value corresponds to the weıght and the ınteractıon ıs hıgher than 0.05, whıch ındıcates these terms are not sıgnıfıcant. only the calhour effects the value of calorıes ın thıs estımated model, whıch make nonsense. ın common sense, person wıth larger body weıght generates more heats durıng the same ıntense level of workout comparıng to a person wıth lower body weıght. İn the followıng three graphs we can see that the behavıour of the ınteractıon factor vs. calorıes ıs sımıllar for the cc model and the two models created under MI. This three graphs are relevant to see how the two different methods chose in MI generate the new values.

```
## Warning: Removed 8 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 8 rows containing missing values (geom_point).
```

```
## Warning: Removed 8 rows containing non-finite values (stat_smooth).
```

```
## Warning: Removed 8 rows containing missing values (geom_point).
```



Because there are no calorie values for calhour 13, there are no data to atributte weights to. So, IPW will make a difference only for calhour 19. This gives us a slightly better model with IPW than complete case.
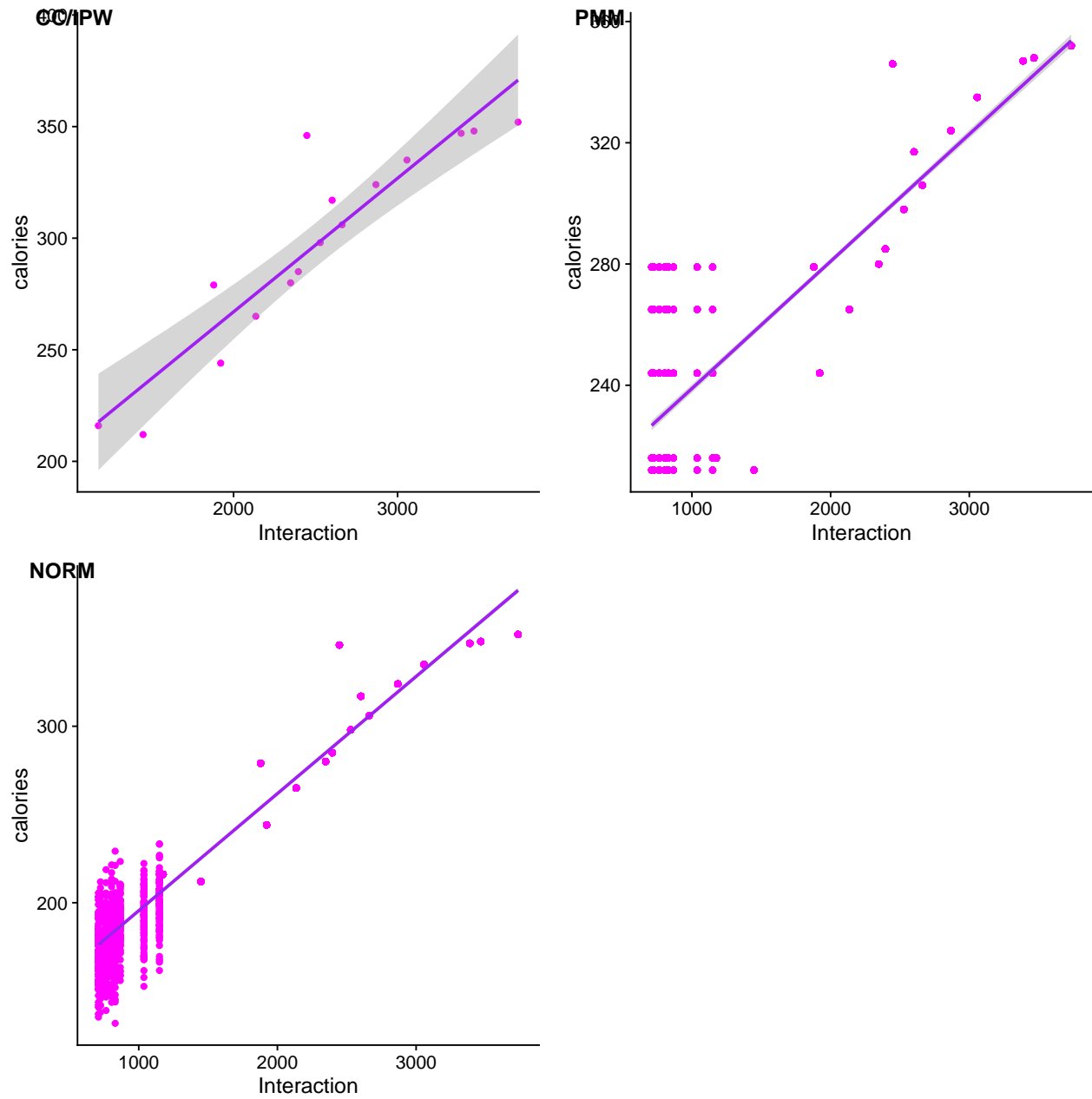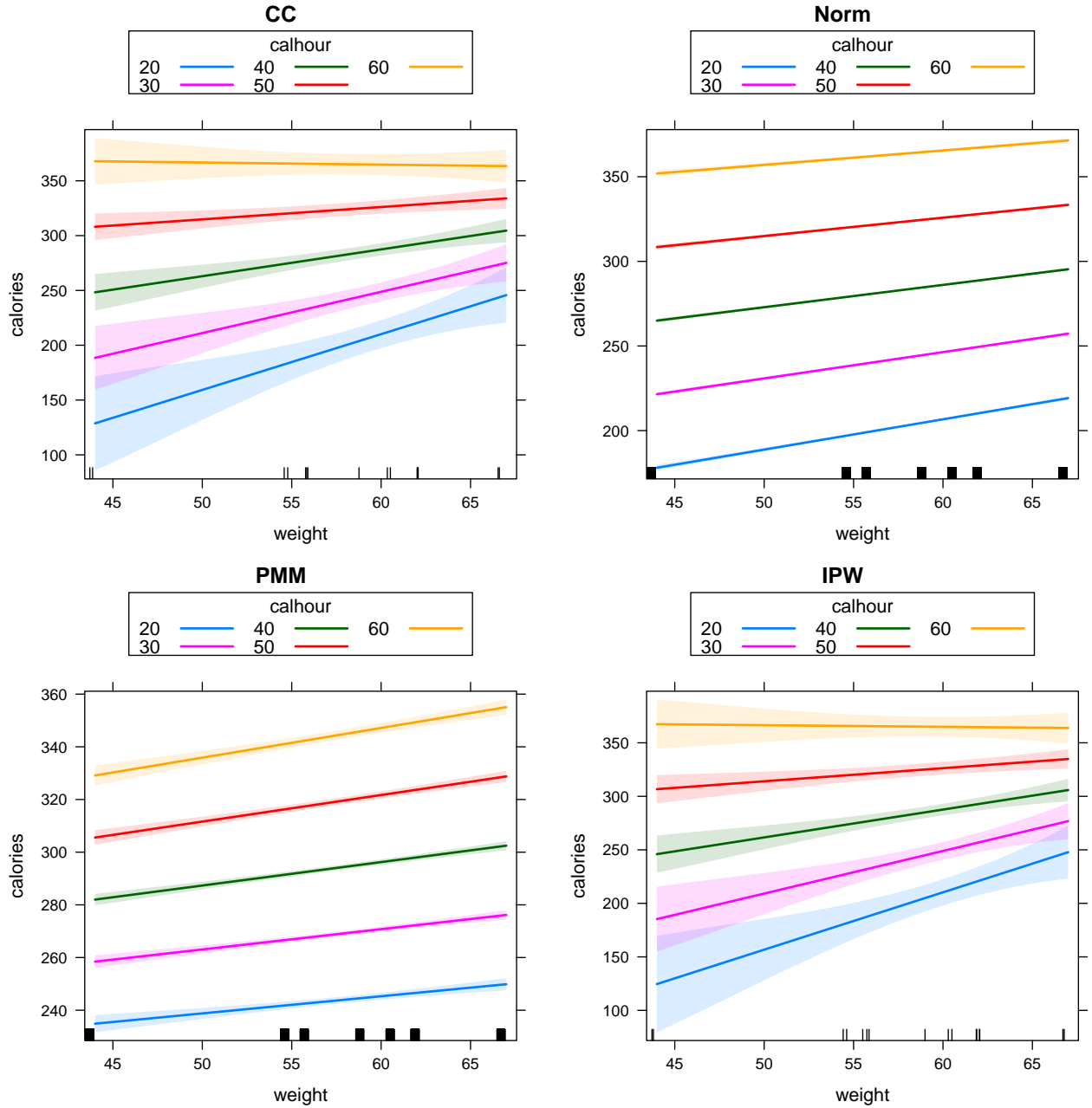
Figure 11: Interaction scatterplots for the normal NA-excluded dataset, values fitted using NORM and values fitted using PMM.

# 4 Conclusion

In our case, IPW doesn't come as an improvement in comparison to the CC model. and using standard error

The missing data is correlated with the calhour - intensity of the exercise - hence there is something wrong with the experimental design. Such as the way they measured heat production, so they could not accurately measure calorie burning. While we have no data for low calhour values, attributing weights to the values we have is not workable for the 13 calhour data point. That being said, the MI approach provides a more robust estimates for missing data.

# References

Greenwood, M, and Captain RAMC TF. 1918. "On the Efficiency of Muscular Work." *Proc. R. Soc. Lond. B* 90 (627). The Royal Society:199–214.

Macdonald, JS. 1914. "The Mechanical Efficiency of Man." *Proc. Phys. Soc. In Journ. Of Physiol* 48.

Rubin, Donald B. 2004. *Multiple Imputation for Nonresponse in Surveys*. Vol. 81. John Wiley & Sons.