

[STAT W4702] Statistical Inference & Modelling Group Project

Babies

12 December 2015

Abstract

Data Set

This project was conducted on the Low Birth Weight dataset collected in 1986 at Baystate Medical Center, Springfield, Massachusetts as a part of a bigger study on the factors influencing newborn infants' health and risk of serious health problems potentially leading to death. This dataset is distributed as a part of MASS library and contains **189 observations** and **10 variables**, among which **bwt** represents the exact amount of newborn infant's weight in grams and is used as the variable of interest we are trying to predict. The other 9 variables stand for different factors related to mothers' physiological parameters, such as age, weight and race, their health-related habits and behavior during pregnancy (smoking habits, presence of uterine irritability and number of physician visits). Also there is a low birth weight indicator **low**, which is defined as a binary variable showing whether the weight of an infant is below 2500 grams or not. Brief description of each variable is provided in the table below.

The goal of our research is to identify relationship between these variables and infant weight and understand the influence of each of them on the explained variable. The project pursue both inferential and predictive goals as it is equally important to be able to obtain inference about factors affecting newborn's health and to be able to react on the potential health risks in a timely manner, when the model predicts the low birth weight outcome for a certain observation. In order to accomplish this goal we tried to fit multiple linear and non-linear models exploring the rationale that could provide the evidence for certain types of models and finding balance between interpretability and predictive power of the model.

Cleaning Dataset

For the purposes of the research the dataset was cleaned in the following way:

- factor variable **race** was assigned with proper labels **white**, **black** and **other**;
- physician visits were converted to a factor variable **ftv** with 3 labels 0, 1 and 2+;
- response is defined as an exact amount of infant's weight from **bwt**;
- all the columns are assigned with meaningful names.

Variable description table and summary statistics of the tidy dataset are provided below.

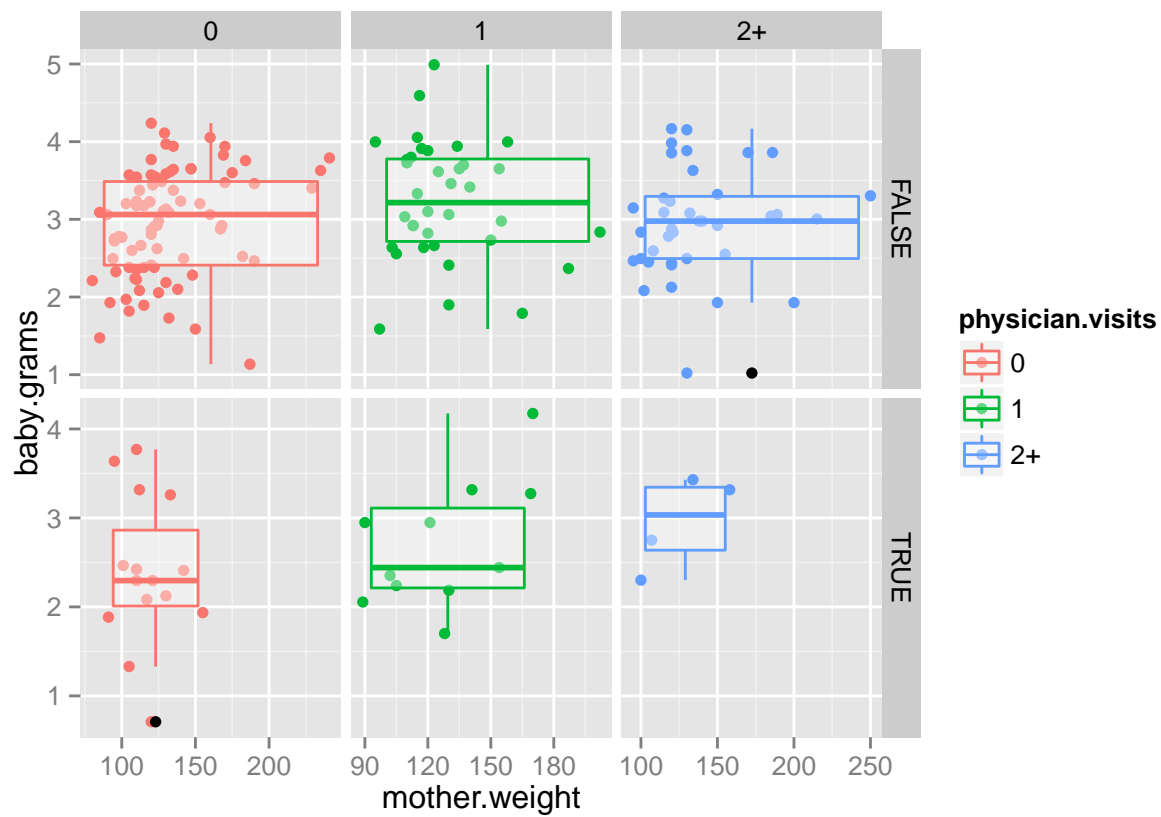
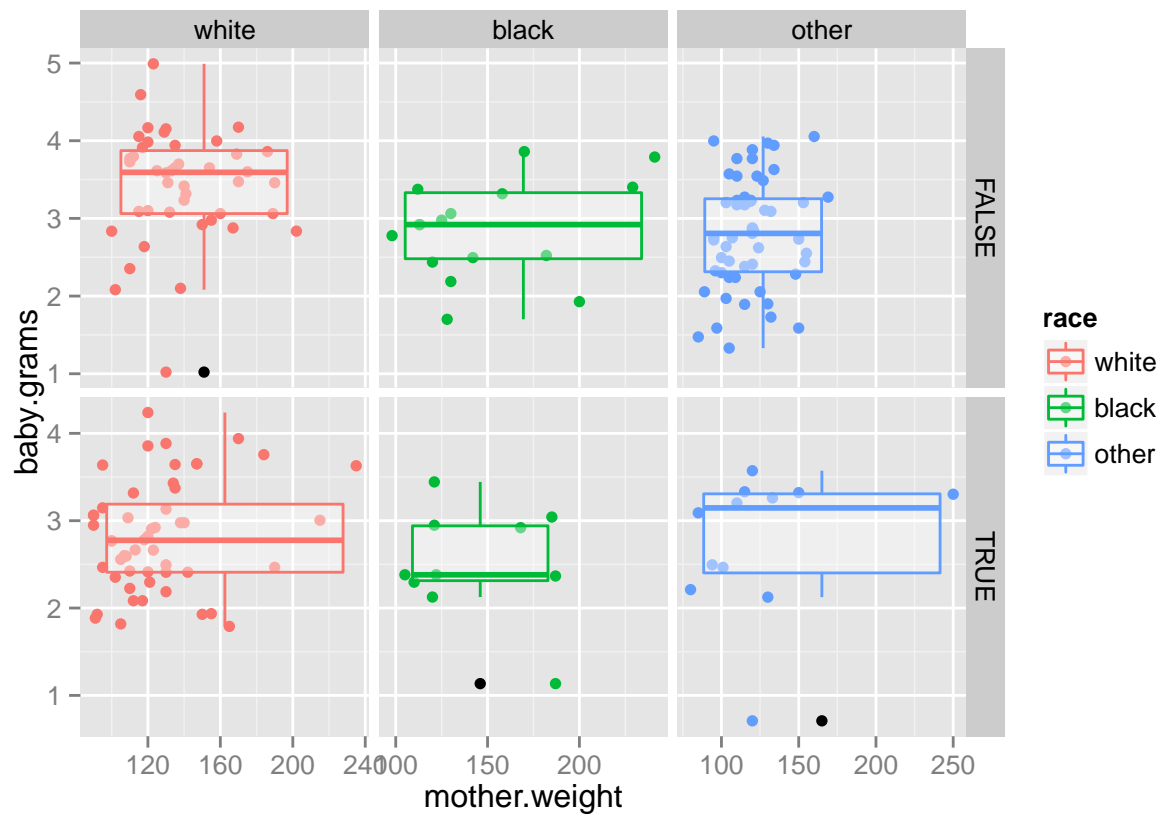
Variable	Description
baby.grams	weight of newborn infant in grams
mother.age	mother's age in years
mother.weight	mother's weight in pounds at last menstrual period
race	mother's race, factor variable with following labels: <i>white</i> , <i>black</i> or <i>other</i>
smoke	smoking status during pregnancy, binary variable
prem.labor	binary variable showing whether mother had premature labors before or not
hypertension	binary variable showing whether mother had hypertension or not

Variable	Description
uterine	binary variable showing presence of uterine irritability
physician.visits	number of physician visits during the first trimester: 0, 1 or 2+

```
##      baby.grams      mother.age      mother.weight      race
##  Min.      :0.709    Min.      :14.00    Min.      : 80.0    white:96
##  1st Qu.:2.414    1st Qu.:19.00    1st Qu.:110.0    black:26
##  Median :2.977    Median :23.00    Median :121.0    other:67
##  Mean   :2.945    Mean   :23.24    Mean   :129.8
##  3rd Qu.:3.487    3rd Qu.:26.00    3rd Qu.:140.0
##  Max.   :4.990    Max.   :45.00    Max.   :250.0
##      smoke      prem.labor      hypertension      uterine
##  Mode :logical  FALSE:159    Mode :logical    Mode :logical
##  FALSE:115      TRUE : 30    FALSE:177        FALSE:161
##  TRUE :74                TRUE :12        TRUE :28
##  NA's :0                NA's :0         NA's :0
##
##
##  physician.visits
##  0 :100
##  1 : 47
##  2+: 42
##
##
##
```

Datset has only 3 quantitative variables, however, as shown in the table below, they do not demonstrate

```
##      baby.grams      mother.age      mother.weight
##  baby.grams      1.00000000  0.09031781      0.1857333
##  mother.age      0.09031781  1.00000000      0.1800732
##  mother.weight  0.18573328  0.18007315      1.0000000
```



```
set.seed(1)
train <- sample(1:nrow(bwt.grams), floor(0.75*nrow(bwt.grams)))
```

```
library(MASS)
data(birthwt)
bwt <- with(birthwt, {
  race <- factor(race, labels = c("white", "black", "other"))
  ptd <- factor(ptl > 0)
  ftv <- factor(ftv)
  levels(ftv)[- (1:2)] <- "2+"
  data.frame(low, age, lwt, race, smoke = (smoke > 0),
             ptd, ht = (ht > 0), ui = (ui > 0), ftv)
})
colnames(bwt) <- c("below.2500", "mother.age",
                  "mother.weight", "race",
                  "smoke", "prem.labor",
                  "hypertension", "uterine",
                  "physician.visits")

bwt.grams <- with(birthwt, {
  bwt <- bwt/1000
  race <- factor(race, labels = c("white", "black", "other"))
  ptd <- factor(ptl > 0)
  ftv <- factor(ftv)
  levels(ftv)[- (1:2)] <- "2+"
  data.frame(bwt, age, lwt, race, smoke = (smoke > 0),
             ptd, ht = (ht > 0), ui = (ui > 0), ftv)
})
colnames(bwt.grams) <- c("baby.grams", "mother.age",
                       "mother.weight", "race",
                       "smoke", "prem.labor",
                       "hypertension", "uterine",
                       "physician.visits")

summary(bwt)
```

```
##    below.2500    mother.age    mother.weight    race
## Min.   :0.0000   Min.   :14.00   Min.    : 80.0   white:96
## 1st Qu.:0.0000   1st Qu.:19.00   1st Qu.:110.0   black:26
## Median :0.0000   Median :23.00   Median :121.0   other:67
## Mean   :0.3122   Mean    :23.24   Mean    :129.8
## 3rd Qu.:1.0000   3rd Qu.:26.00   3rd Qu.:140.0
## Max.    :1.0000   Max.    :45.00   Max.    :250.0
##    smoke      prem.labor  hypertension    uterine
## Mode :logical FALSE:159   Mode :logical  Mode :logical
## FALSE:115    TRUE : 30   FALSE:177     FALSE:161
## TRUE :74      NA's :0    TRUE :12      TRUE :28
## NA's :0      NA's :0    NA's :0       NA's :0
##
##
## physician.visits
## 0 :100
## 1 : 47
## 2+: 42
```

```
##  
##  
##
```

```
summary(bwt.grams)
```

```
##      baby.grams      mother.age      mother.weight      race  
## Min.   :0.709      Min.   :14.00      Min.   : 80.0      white:96  
## 1st Qu.:2.414      1st Qu.:19.00      1st Qu.:110.0      black:26  
## Median :2.977      Median :23.00      Median :121.0      other:67  
## Mean   :2.945      Mean   :23.24      Mean   :129.8  
## 3rd Qu.:3.487      3rd Qu.:26.00      3rd Qu.:140.0  
## Max.   :4.990      Max.   :45.00      Max.   :250.0  
##      smoke      prem.labor      hypertension      uterine  
## Mode :logical      FALSE:159      Mode :logical      Mode :logical  
## FALSE:115      TRUE : 30      FALSE:177      FALSE:161  
## TRUE :74      TRUE :12      TRUE :28  
## NA's :0      NA's :0      NA's :0  
##  
##  
## physician.visits  
## 0 :100  
## 1 : 47  
## 2+: 42  
##  
##  
##
```

```
bwt[0:10,]
```

```
##      below.2500      mother.age      mother.weight      race      smoke      prem.labor      hypertension  
## 1              0              19              182      black      FALSE              FALSE              FALSE  
## 2              0              33              155      other      FALSE              FALSE              FALSE  
## 3              0              20              105      white      TRUE               FALSE              FALSE  
## 4              0              21              108      white      TRUE               FALSE              FALSE  
## 5              0              18              107      white      TRUE               FALSE              FALSE  
## 6              0              21              124      other      FALSE              FALSE              FALSE  
## 7              0              22              118      white      FALSE              FALSE              FALSE  
## 8              0              17              103      other      FALSE              FALSE              FALSE  
## 9              0              29              123      white      TRUE               FALSE              FALSE  
## 10             0              26              113      white      TRUE               FALSE              FALSE  
##      uterine      physician.visits  
## 1      TRUE              0  
## 2     FALSE              2+  
## 3     FALSE              1  
## 4      TRUE              2+  
## 5      TRUE              0  
## 6     FALSE              0  
## 7     FALSE              1  
## 8     FALSE              1  
## 9     FALSE              1  
## 10    FALSE              0
```

```
bwt.grams[0:10,]
```

```
##      baby.grams mother.age mother.weight  race smoke prem.labor hypertension
## 1      2.523      19      182 black FALSE      FALSE      FALSE
## 2      2.551      33      155 other FALSE      FALSE      FALSE
## 3      2.557      20      105 white  TRUE      FALSE      FALSE
## 4      2.594      21      108 white  TRUE      FALSE      FALSE
## 5      2.600      18      107 white  TRUE      FALSE      FALSE
## 6      2.622      21      124 other FALSE      FALSE      FALSE
## 7      2.637      22      118 white FALSE      FALSE      FALSE
## 8      2.637      17      103 other FALSE      FALSE      FALSE
## 9      2.663      29      123 white  TRUE      FALSE      FALSE
## 10     2.665      26      113 white  TRUE      FALSE      FALSE
##      uterine physician.visits
## 1      TRUE      0
## 2     FALSE      2+
## 3     FALSE      1
## 4      TRUE      2+
## 5      TRUE      0
## 6     FALSE      0
## 7     FALSE      1
## 8     FALSE      1
## 9     FALSE      1
## 10    FALSE      0
```

```
attach(bwt.grams)
```

```
library(leaps)
regfit.full=regsubsets(baby.grams~., bwt.grams, nvmax =19)
reg.summary = summary(regfit.full)
reg.summary$rsq
```

```
## [1] 0.08061477 0.11225032 0.14782772 0.18905712 0.21364404 0.24039446
## [7] 0.25042689 0.25537670 0.25647316 0.25682243
```

```
par(mfrow =c(2,2))
plot(reg.summary$rsq ,xlab=" Number of Variables ",ylab=" RSS", type="l")
plot(reg.summary$adjr2 ,xlab=" Number of Variables ", ylab=" Adjusted RSq",type="l")
max.adj2=which.max (reg.summary$adjr2)
max.adj2
```

```
## [1] 8
```

```
points (max.adj2, reg.summary$adjr2[max.adj2], col ="red",cex =2, pch =20)
```

```
plot(reg.summary$cp ,xlab=" Number of Variables ", ylab="Cp", type='l')
min.cp= which.min (reg.summary$cp )
min.cp
```

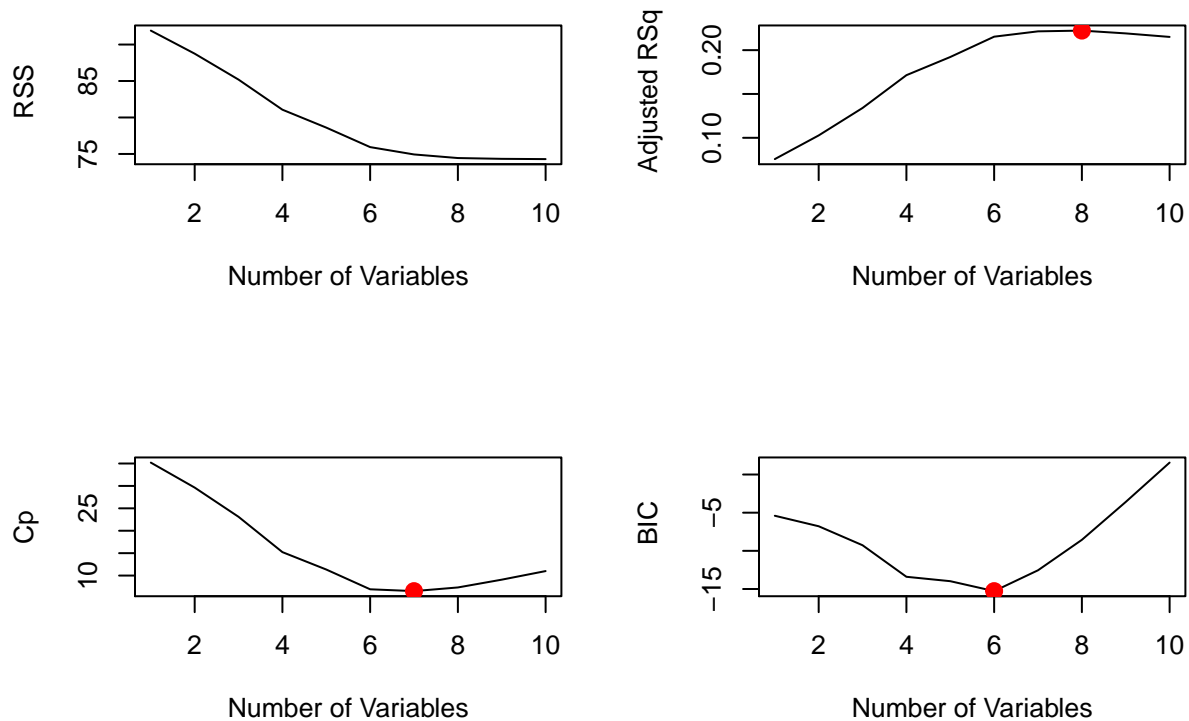
```
## [1] 7
```

```
points (min.cp, reg.summary$cp[min.cp], col ="red",cex =2, pch =20)
```

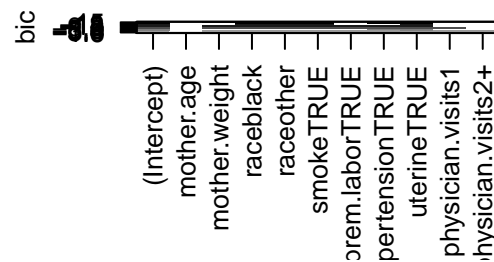
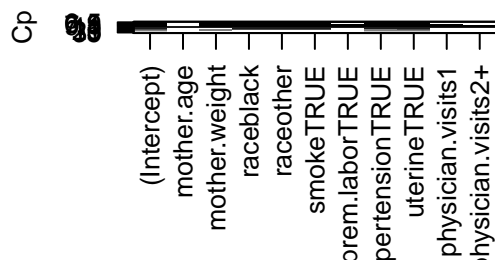
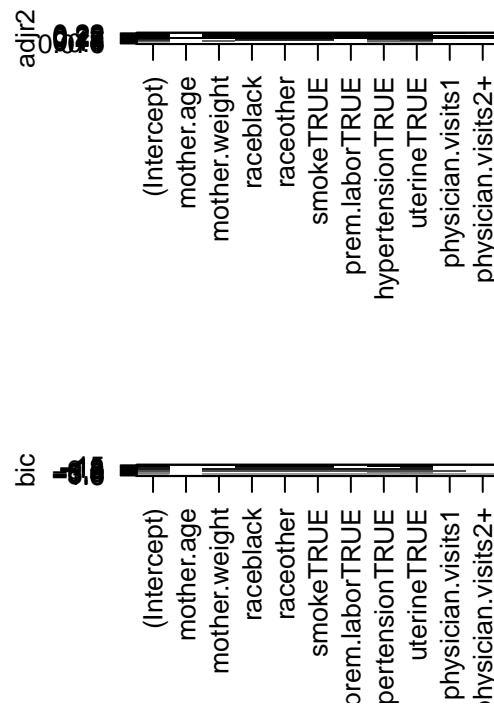
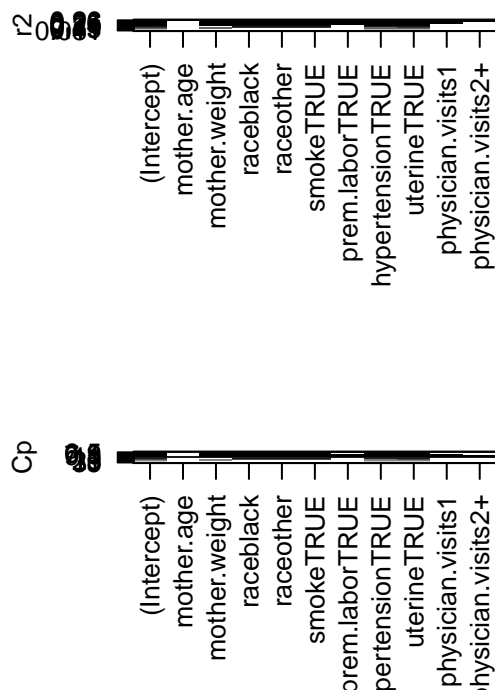
```
min.bic = which.min(reg.summary$bic)
min.bic
```

```
## [1] 6
```

```
plot(reg.summary$bic ,xlab=" Number of Variables ",ylab=" BIC", type='l')
points (min.bic, reg.summary$bic [min.bic], col =" red",cex =2, pch =20)
```



```
plot(regfit.full ,scale ="r2")
plot(regfit.full ,scale ="adjr2")
plot(regfit.full ,scale ="Cp")
plot(regfit.full ,scale ="bic")
```



```
coef(regfit.full, max.adjR2)
```

```
##      (Intercept)      mother.weight      raceblack      raceother
##      2.799714010      0.004194539      -0.453359173      -0.305169792
##      smokeTRUE      prem.laborTRUE      hypertensionTRUE      uterineTRUE
##      -0.294468372      -0.235263456      -0.577857003      -0.478599299
## physician.visits1
##      0.125220667
```

```
coef(regfit.full, min.cp)
```

```
##      (Intercept)      mother.weight      raceblack      raceother
##      2.871512227      0.004043831      -0.465601219      -0.333878191
##      smokeTRUE      prem.laborTRUE      hypertensionTRUE      uterineTRUE
##      -0.325081991      -0.207834528      -0.573799253      -0.491143889
```

```
coef(regfit.full, min.bic)
```

```
##      (Intercept)      mother.weight      raceblack      raceother
##      2.83726392      0.00424155      -0.47505760      -0.34815038
##      smokeTRUE      hypertensionTRUE      uterineTRUE
##      -0.35632095      -0.58519312      -0.52552390
```

```
classfit.full=regsubsets(below.2500~., bwt, nvmax =19)
class.summary = summary(classfit.full)
class.summary$rsq
```

```
## [1] 0.07279919 0.09555397 0.12812223 0.14603025 0.16130952 0.17290333
## [7] 0.18432185 0.19036572 0.19240164 0.19259390
```



```

par(mfrow =c(2,2))
plot(class.summary$rss ,xlab=" Number of Variables ",ylab=" RSS", type="l")
plot(class.summary$adjr2 ,xlab = " Number of Variables ", ylab=" Adjusted RSq",type="l")
max.adj2=which.max (class.summary$adjr2)
max.adj2

```

```
## [1] 8
```

```
points (max.adj2, class.summary$adjr2[max.adj2], col ="red",cex =2, pch =20)
```

```

plot(class.summary$cp ,xlab = " Number of Variables ", ylab="Cp", type='l')
min.cp= which.min (class.summary$cp )
min.cp

```

```
## [1] 7
```

```
points (min.cp, class.summary$cp[min.cp], col ="red",cex =2, pch =20)
```

```

min.bic = which.min(class.summary$bic)
min.bic

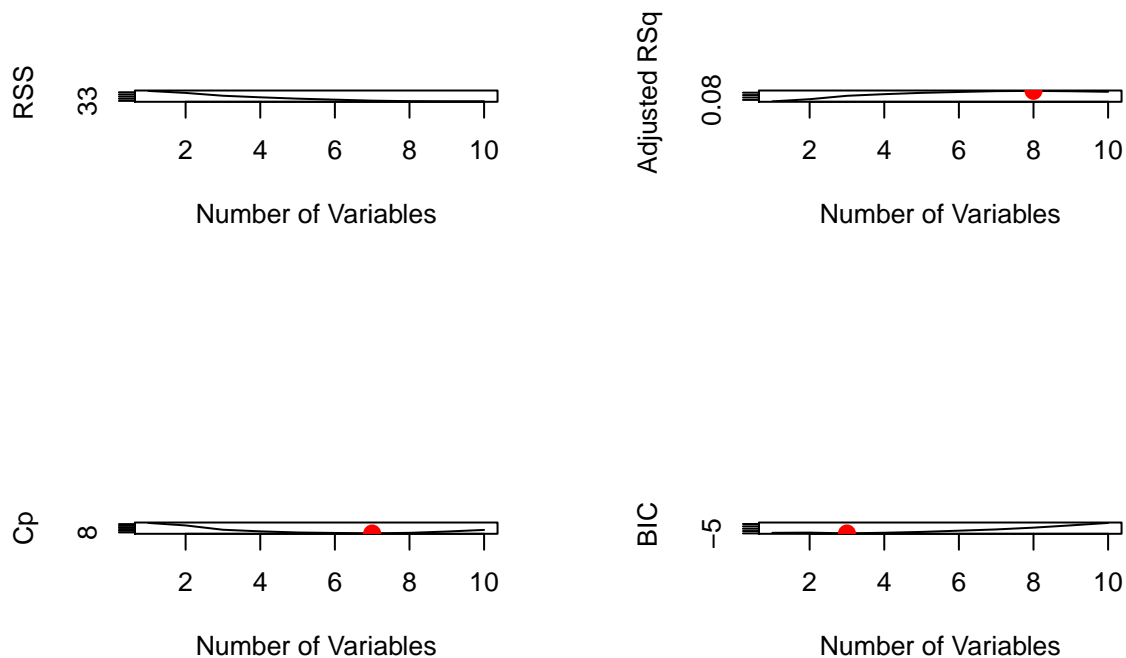
```

```
## [1] 3
```

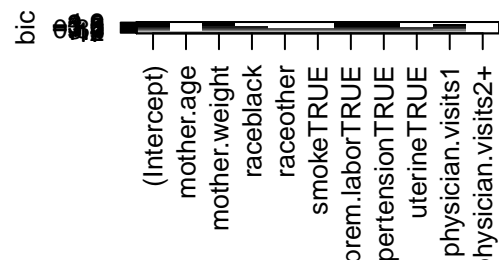
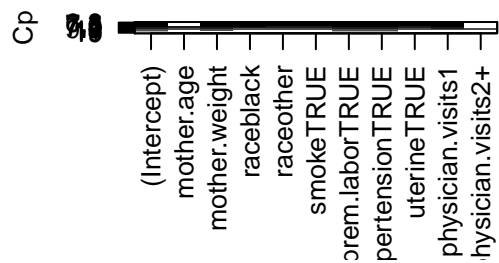
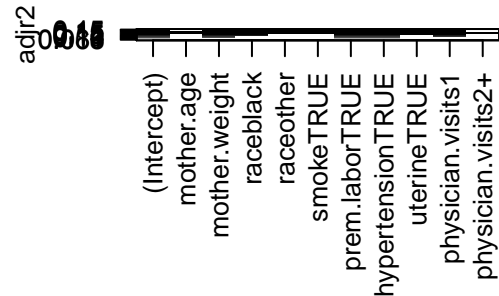
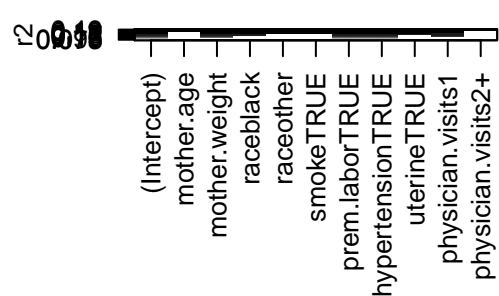
```

plot(class.summary$bic ,xlab=" Number of Variables ",ylab=" BIC", type='l')
points (min.bic, class.summary$bic [min.bic], col =" red",cex =2, pch =20)

```



```
plot(classfit.full, scale = "r2")
plot(classfit.full, scale = "adjr2")
plot(classfit.full, scale = "Cp")
plot(classfit.full, scale = "bic")
```



```
coef(classfit.full, max.adjR2)
```

```
##      (Intercept)      mother.weight      raceblack      raceother
##      0.47731870      -0.00269524      0.21446267      0.11814439
##      smokeTRUE      prem.laborTRUE      hypertensionTRUE      uterineTRUE
##      0.12582999      0.26509425      0.36294635      0.14095381
## physician.visits1
##      -0.08816014
```

```
coef(classfit.full, min.cp)
```

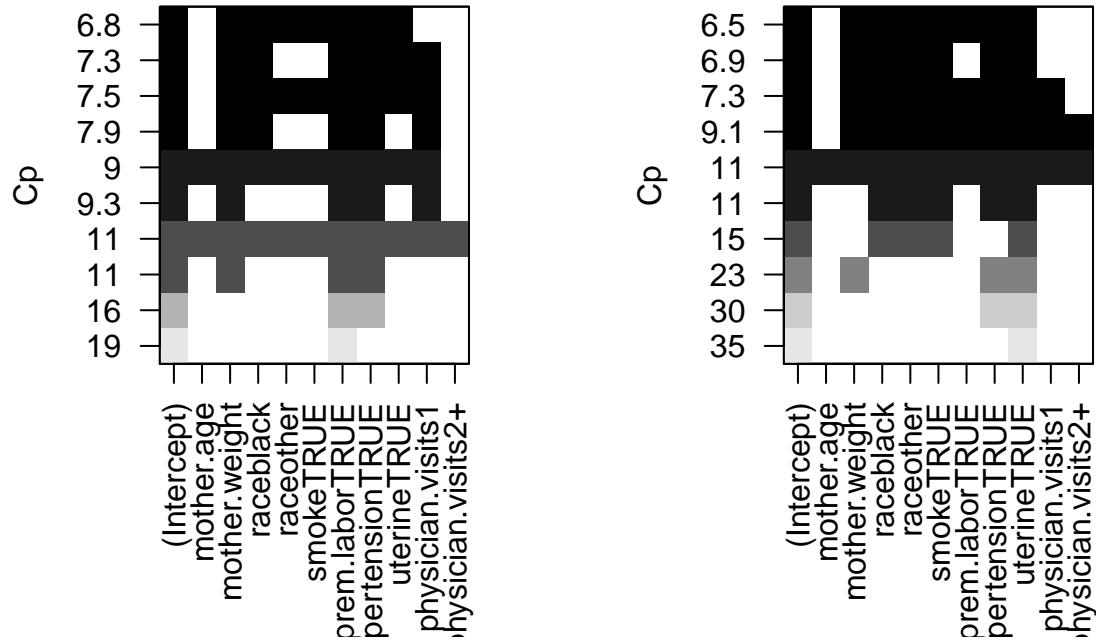
```
##      (Intercept)      mother.weight      raceblack      raceother
##      0.42677003      -0.002589136      0.223081542      0.138356202
##      smokeTRUE      prem.laborTRUE      hypertensionTRUE      uterineTRUE
##      0.147383151      0.245783235      0.360089535      0.149785677
```

```
coef(classfit.full, min.bic)
```

```
##      (Intercept)      mother.weight      prem.laborTRUE      hypertensionTRUE
##      0.607928770      -0.002842709      0.313205471      0.370930320
```

```
#Let's compare classification and regression
```

```
par(mfrow =c(1,2))
plot(classfit.full ,scale ="Cp")
plot(regfit.full ,scale ="Cp")
```



```
#Logistic regression with the predictors selected by best subset
```

```
log.fit = glm( below.2500~ mother.weight+race+smoke+hypertension+uterine, family = binomial, data=bwt[t,])
summary(log.fit)
```

```
##
## Call:
## glm(formula = below.2500 ~ mother.weight + race + smoke + hypertension +
##      uterine, family = binomial, data = bwt[train, ])
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.8144  -0.7984  -0.4335   0.8262   2.1800
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -0.143646   1.153607  -0.125  0.90090
## mother.weight -0.019410   0.008422  -2.305  0.02119 *
## raceblack     1.671922   0.653759   2.557  0.01055 *
## raceother     1.395900   0.561766   2.485  0.01296 *
## smokeTRUE     1.543006   0.511265   3.018  0.00254 **
## hypertensionTRUE 2.023435   0.777850   2.601  0.00929 **
## uterineTRUE    1.041207   0.545658   1.908  0.05637 .
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
```

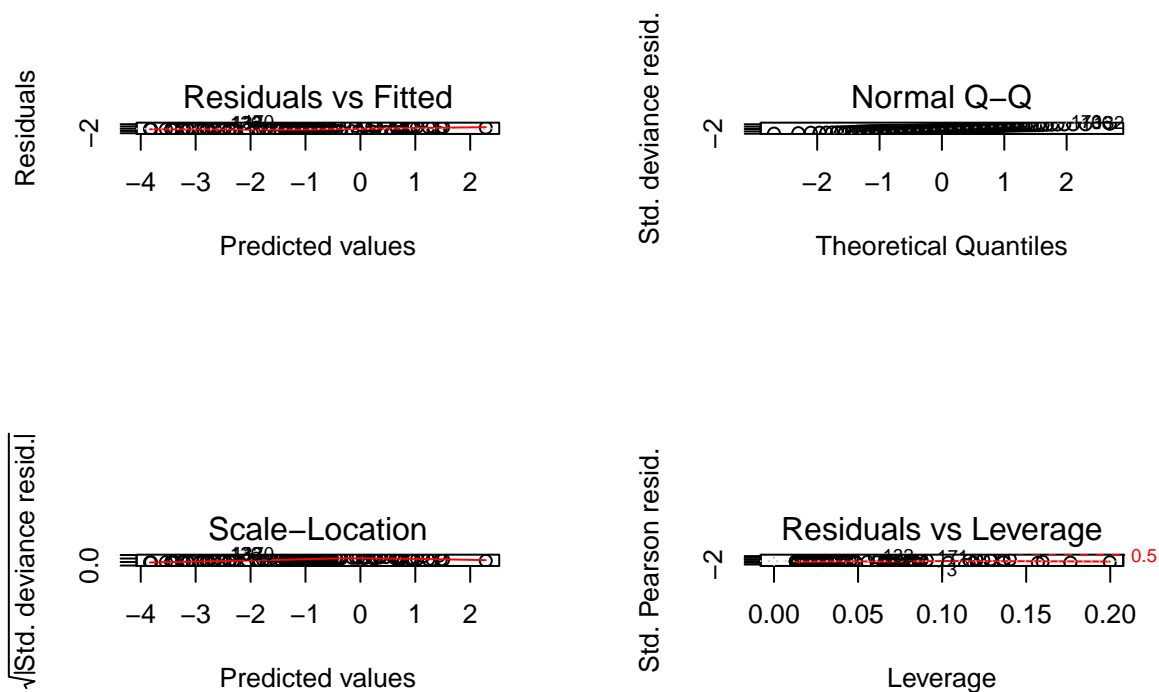
```
## Null deviance: 175.05 on 140 degrees of freedom
## Residual deviance: 142.72 on 134 degrees of freedom
## AIC: 156.72
##
## Number of Fisher Scoring iterations: 5
```

```
confint(log.fit)
```

```
## Waiting for profiling to be done...
```

```
##              2.5 %      97.5 %
## (Intercept) -2.32618747 2.220864159
## mother.weight -0.03731747 -0.004067736
## raceblack    0.41072752 3.005448309
## raceother    0.33601380 2.561460286
## smokeTRUE    0.58304495 2.610272893
## hypertensionTRUE 0.54469293 3.659423231
## uterineTRUE  -0.02878761 2.132040177
```

```
par(mfrow = c(2, 2))
plot(log.fit)
```



```
pred.train <- predict(log.fit, type = "response")
low.train <- sapply(pred.train, function(x) {ifelse(x > 0.5, 1, 0)})
table(low.train, bwt$below.2500[train])
```

```
##
## low.train  0  1
##           0 87 25
##           1 10 19
```

```
mean(low.train == bwt$below.2500[train])
```

```
## [1] 0.751773
```

```
pred.test <- predict(log.fit, newdata = bwt[-train, -1], type = "response")
low.test <- sapply(pred.test, function(x) {ifelse(x > 0.5, 1, 0)})
table(low.test, bwt$below.2500[-train])
```

```
##
## low.test  0  1
##           0 29 11
##           1  4  4
```

```
mean(low.test == bwt$below.2500[-train])
```

```
## [1] 0.6875
```

```
#Linear regression with the predictors selected by best subset
```

```
lm.fit = lm( baby.grams~ mother.weight+race+smoke+hypertension+uterine, data=bwt.grams)
summary(lm.fit)
```

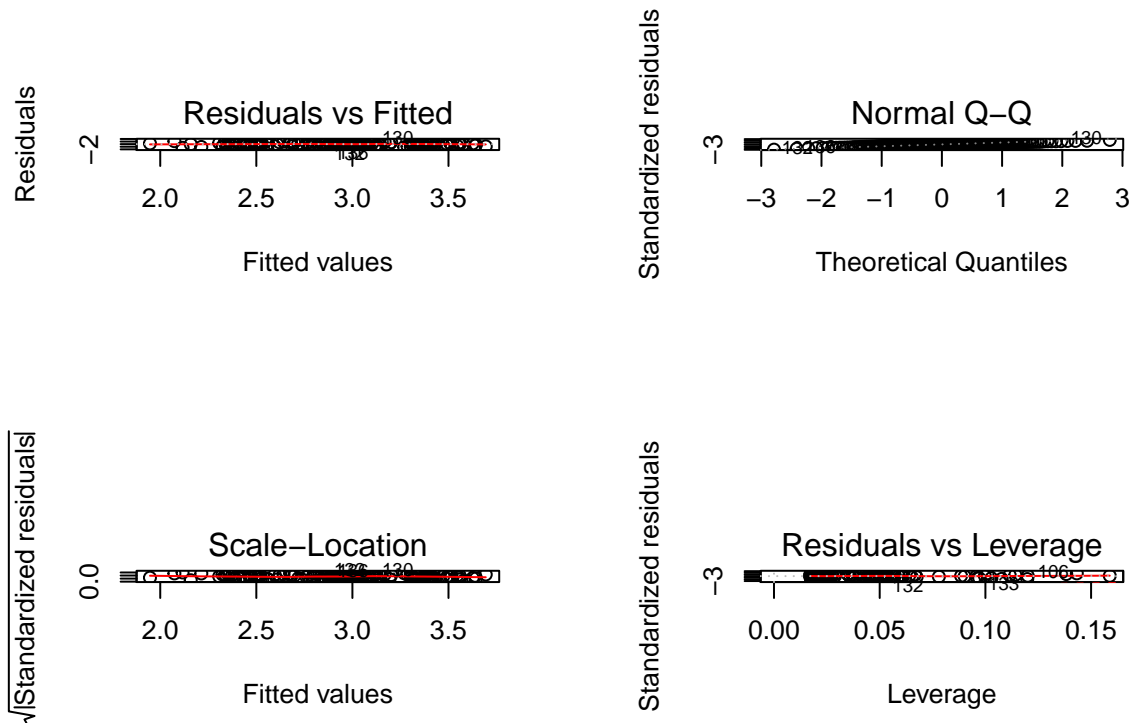
```
##
## Call:
## lm(formula = baby.grams ~ mother.weight + race + smoke + hypertension +
##      uterine, data = bwt.grams)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.84214 -0.43319  0.06709  0.45921  1.63103
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    2.837264    0.243676   11.644 < 2e-16 ***
## mother.weight    0.004242    0.001675    2.532 0.012198 *
## raceblack      -0.475058    0.145603   -3.263 0.001318 **
## raceother      -0.348150    0.112361   -3.099 0.002254 **
## smokeTRUE      -0.356321    0.103444   -3.445 0.000710 ***
## hypertensionTRUE -0.585193    0.199644   -2.931 0.003810 **
## uterineTRUE     -0.525524    0.134675   -3.902 0.000134 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.6459 on 182 degrees of freedom
## Multiple R-squared:  0.2404, Adjusted R-squared:  0.2154
## F-statistic:  9.6 on 6 and 182 DF, p-value: 3.601e-09
```

```
confint(lm.fit)
```

```
##              2.5 %      97.5 %
## (Intercept)  2.3564706569  3.318057183
```

```
## mother.weight      0.0009358509  0.007547249
## raceblack         -0.7623440159 -0.187771193
## raceother         -0.5698476393 -0.126453123
## smokeTRUE         -0.5604237850 -0.152218115
## hypertensionTRUE -0.9791080814 -0.191278160
## uterineTRUE        -0.7912496587 -0.259798136
```

```
par(mfrow = c(2, 2))
plot(lm.fit)
```

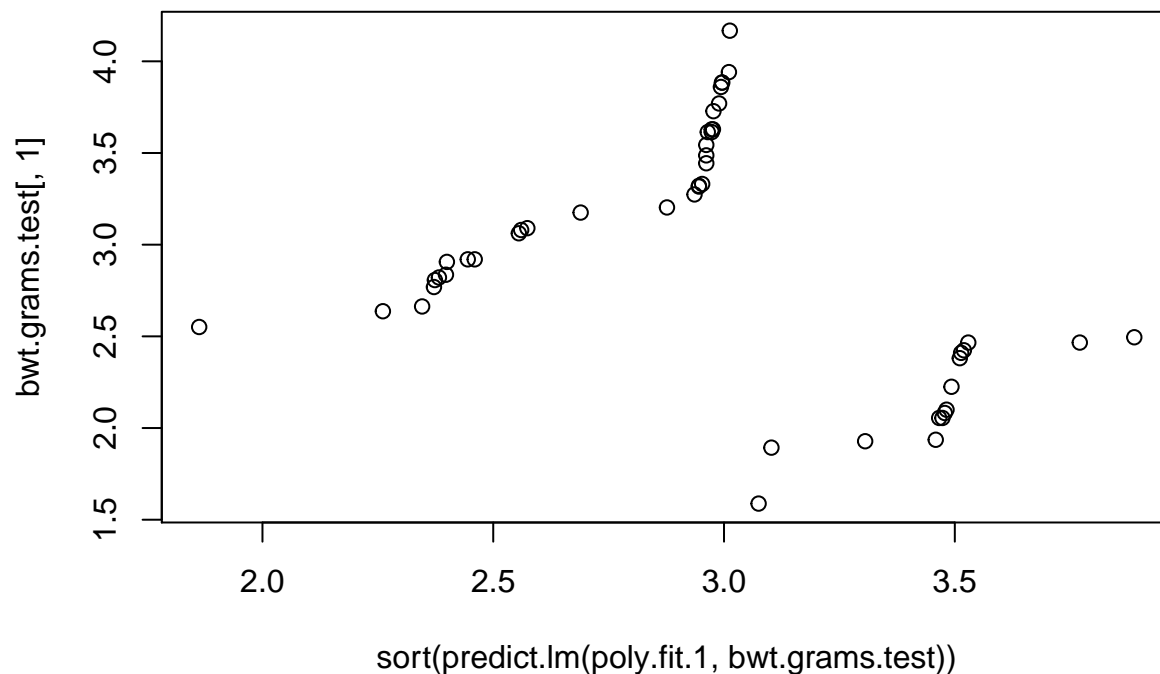


```
#Create train and test
set.seed(1)
train <- sample(1:nrow(bwt.grams), floor(0.75*nrow(bwt.grams)))
bwt.grams.train <- bwt.grams[train,]
bwt.grams.test <- bwt.grams[-train,]

#Polynomial fit for best subset
poly.fit.1 = lm(baby.grams ~ hypertension + uterine + smoke + race + poly(mother.weight, 2), data = bwt.grams.train)
mean((predict.lm(poly.fit.1, bwt.grams.test) - bwt.grams.test[,1])^2)
```

```
## [1] 0.4813745
```

```
plot(sort(predict.lm(poly.fit.1, bwt.grams.test)), bwt.grams.test[,1])
```



```
anova(poly.fit.1)
```

```
## Analysis of Variance Table
##
## Response: baby.grams
##
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
hypertension	1	2.827	2.8270	6.9365	0.009446 **
uterine	1	8.145	8.1454	19.9863	1.654e-05 ***
smoke	1	3.026	3.0264	7.4259	0.007294 **
race	2	9.074	4.5371	11.1328	3.386e-05 ***
poly(mother.weight, 2)	2	2.363	1.1813	2.8986	0.058590 .
Residuals	133	54.204	0.4075		

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

When we fit a polynomial model on the predictors obtained from best subset, we observe a Mean Squared Error of 0.4813745. The smaller the Mean Squared Error, the closer the fit is to the data. But, as the value of MSE is high, it suggests that this model does not provide a good fit for the data. The plot also shows that there are irregularities in the prediction and that the polynomial model of degree 2 obtained by using predictors suggested by the best subset is not sufficient. When we perform Analysis of Variance (ANOVA) on the polynomial fit, we see that, the *p-values* for all the predictors - except `mother.weight` - are less than 0.05 and thus, the NULL hypothesis that these variables affect the baby weight at birth can be rejected.

Different models were tried by increasing the degree of the polynomial but still using the predictors suggested by the best subset and the following results were obtained:

```
poly.fit.2 = lm(baby.grams ~ hypertension + uterine + smoke + race + poly(mother.weight, 3), data = bwt)
mean((predict.lm(poly.fit.2, bwt.grams.test) - bwt.grams.test[,1])^2)
```

```
## [1] 0.4640868
```

```
poly.fit.3 = lm(baby.grams ~ hypertension + uterine + smoke + race + poly(mother.weight, 4), data = bwt)
mean((predict.lm(poly.fit.3, bwt.grams.test) - bwt.grams.test[,1])^2)
```

```
## [1] 0.4619314
```

```
anova(poly.fit.1, poly.fit.2, poly.fit.3)
```

```
## Analysis of Variance Table
##
## Model 1: baby.grams ~ hypertension + uterine + smoke + race + poly(mother.weight,
##      2)
## Model 2: baby.grams ~ hypertension + uterine + smoke + race + poly(mother.weight,
##      3)
## Model 3: baby.grams ~ hypertension + uterine + smoke + race + poly(mother.weight,
##      4)
##      Res.Df    RSS Df Sum of Sq    F Pr(>F)
## 1      133 54.204
## 2      132 53.825  1   0.37915 0.9239 0.3382
## 3      131 53.761  1   0.06393 0.1558 0.6937
```

We note that as the degree of the polynomial increases, the MSE decreases, but the drop is not significant, suggesting that these predictors are not sufficient enough to predict the correct baby weight. Performing the ANOVA test to compare how the three models perform with respect to each other, we observe high *p-values* which state that the none of the models are good enough.

When we remove the predictors with very low *p-values*, which were suggested by the best subset - namely **smoke**, **race** and add other predictors which were rejected by the best-subset, namely - **mother.age**, **prem.labor** and **physician.visits**, we see that the Mean Squared Error starts to decrease. A low MSE denotes a better fit. Thus, the predictors which were rejected by the best subset selection, were actually significant in predicting the correct birthweight.

```
poly.fit.4 = lm(baby.grams ~ hypertension + uterine + poly(mother.age,2) + poly(mother.weight,3), data = bwt)
mean((predict.lm(poly.fit.4, bwt.grams.test) - bwt.grams.test[,1])^2)
```

```
## [1] 0.3890751
```

```
poly.fit.5 = lm(baby.grams ~ hypertension + uterine + smoke + prem.labor + poly(mother.age,2) + poly(mother.weight,3), data = bwt)
mean((predict.lm(poly.fit.5, bwt.grams.test) - bwt.grams.test[,1])^2)
```

```
## [1] 0.3865828
```

```
poly.fit.6 = lm(baby.grams ~ hypertension + uterine + smoke + prem.labor + poly(mother.age,2) + poly(mother.weight,3) + physician.visits, data = bwt)
mean((predict.lm(poly.fit.6, bwt.grams.test) - bwt.grams.test[,1])^2)
```

```
## [1] 0.3214657
```