

TMA4315: Project 1

Jim Totland, Martin Gudahl Tufte

9/8/2021

Problem 1

a)

Since the response variables $y_i \sim \text{Bernoulli}(\pi_i)$, where $\pi_i = \Pr(y_i = 1 \mid \mathbf{x}_i)$. The conditional mean is given by $Ey_i = \pi_i$, which is connected to the covariates via the following relationship:

$$\mathbf{x}_i^T \boldsymbol{\beta} =: \eta_i = \Phi^{-1}(\pi_i),$$

or equivalently: $\pi_i = \Phi(\eta_i)$. This results in the likelihood function

$$\begin{aligned} L(\boldsymbol{\beta}) &= \prod_{i=1}^n \pi_i^{y_i} (1 - \pi_i)^{1-y_i} \\ &= \prod_{i=1}^n \Phi(\eta_i)^{y_i} (1 - \Phi(\eta_i))^{1-y_i}. \end{aligned}$$

Thus, the log-likelihood becomes

$$l(\boldsymbol{\beta}) := \ln(L(\boldsymbol{\beta})) = \sum_{i=1}^n \underbrace{y_i \ln(\Phi(\eta_i)) + (1 - y_i) \ln(1 - \Phi(\eta_i))}_{=l_i(\boldsymbol{\beta})} = \sum_{i=1}^n l_i(\boldsymbol{\beta}).$$

To find the score function, we calculate

$$\begin{aligned} \frac{\partial l_i(\boldsymbol{\beta})}{\partial \boldsymbol{\beta}} &= \frac{y_i}{\Phi(\eta_i)} \frac{\partial \Phi(\eta_i)}{\partial \boldsymbol{\beta}} - \frac{1 - y_i}{1 - \Phi(\eta_i)} \frac{\partial \Phi(\eta_i)}{\partial \boldsymbol{\beta}} \\ &= \frac{y_i}{\Phi(\eta_i)} \phi(\eta_i) \mathbf{x}_i - \frac{1 - y_i}{1 - \Phi(\eta_i)} \phi(\eta_i) \mathbf{x}_i \\ &= \frac{y_i(1 - \Phi(\eta_i)) - (1 - y_i)\Phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \phi(\eta_i) \mathbf{x}_i \\ &= \frac{y_i - \Phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \phi(\eta_i) \mathbf{x}_i. \end{aligned}$$

Consequently, the score function is given by

$$\mathbf{s}(\boldsymbol{\beta}) = \sum_{i=1}^n \frac{y_i - \Phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \phi(\eta_i) \mathbf{x}_i.$$

Next, we find the expected Fisher information, $F(\boldsymbol{\beta})$. We find it by using the result

$$\begin{aligned}
F(\beta) &= \text{Var}(\mathbf{s}(\beta)) = \text{Var}\left(\sum_{i=1}^n \frac{y_i - \Phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \phi(\eta_i) \mathbf{x}_i\right) \\
&= \sum_{i=1}^n \left[\frac{\phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \right]^2 \text{Var}(y_i \mathbf{x}_i) = \sum_{i=1}^n \left[\frac{\phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \right]^2 \mathbf{x}_i \text{Var}(y_i) \mathbf{x}_i^T \\
&= \sum_{i=1}^n \left[\frac{\phi(\eta_i)}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \right]^2 \pi_i(1 - \pi_i) \mathbf{x}_i \mathbf{x}_i^T = \sum_{i=1}^n \frac{\phi(\eta_i)^2}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \mathbf{x}_i \mathbf{x}_i^T,
\end{aligned}$$

Where in the third equality we have used that the y_i 's are independent. The expected Fisher information can also be verified to have this expression by the relationship

$$F(\beta) = \sum_{i=1}^n \frac{h'(\eta_i)^2}{\text{Var}(y_i)} \mathbf{x}_i \mathbf{x}_i^T,$$

where $h'(\eta_i) = \Phi'(\eta_i) = \phi(\eta_i)$ and $\text{Var}(y_i) = \pi_i(1 - \pi_i) = \Phi(\eta_i)(1 - \Phi(\eta_i))$.

b)

The expected Fisher information is given by

$$F(\beta) = \sum_{i=1}^n \frac{\phi(\eta_i)^2}{\Phi(\eta_i)(1 - \Phi(\eta_i))} \mathbf{x}_i \mathbf{x}_i^T = \mathbf{x}^T W \mathbf{x},$$

where $W = \text{diag}\left(\frac{\phi(\eta_i)^2}{\Phi(\eta_i)(1 - \Phi(\eta_i))}\right)$.

The Fisher scoring algorithm states that the next iterate is given by

$$\beta^{(t+1)} = \beta^{(t)} + F(\beta^{(t)})^{-1} \mathbf{s}(\beta^{(t)}).$$

Inserting the expected Fisher information and the score function we get

$$\beta^{(t+1)} = (\mathbf{x}^T W^{(t)} \mathbf{x})^{-1} \mathbf{x}^T W^{(t)} \tilde{\mathbf{y}}^{(t)},$$

where the working response vector $\tilde{\mathbf{y}}^{(t)}$ has element i given by

$$\tilde{y}_i^{(t)} = \mathbf{x}_i^T \beta^{(t)} + \frac{y_i - h(\mathbf{x}_i^T \beta^{(t)})}{h'(\mathbf{x}_i^T \beta^{(t)})} = \eta_i^{(t)} + \frac{y_i - \Phi(\eta_i^{(t)})}{\phi(\eta_i^{(t)})}.$$

The deviance is defined as

$$D = -2 l(\hat{\beta}) + 2 l(\text{saturated model}).$$

Implementing myglm in R:

```

Phi <- function(x) return (pnorm(x))
phi <- function(x) return (dnorm(x))

myglm <- function(formula, data, start = NULL){
  # response variable

```

```

resp <- all.vars(formula)[1]
y <- as.matrix( data[resp] )

# model matrix
X <- model.matrix(formula, data)
n <- dim(X)[1]
p <- dim(X)[2]

# starting beta
if (is.null(start)){
  beta = rep(0, p)
}
else {
  beta = start
}

# Fisher scoring algorithm
max_iter <- 50
tol <- 1e-10
iter <- 0
rel.err <- Inf

while (rel.err > tol & iter < max_iter){
  # calculate eta, y tilde, W
  eta <- X %%% beta
  y.tilde <- eta + (y - Phi(eta)) / (phi(eta))
  W <- diag( as.vector(phi(eta)^2 / (Phi(eta)*(1-Phi(eta)))), n, n )

  # update beta
  A <- t(X) %%% W %%% X
  b <- t(X) %%% W %%% y.tilde
  beta.new <- solve(A, b)

  iter <- iter + 1
  rel.err <- max(abs(beta.new - beta) / abs(beta.new))
  beta <- beta.new
}

Estimate <- beta
F.inv <- solve(A)
Std.Error <- sqrt(diag(F.inv))

return (list("coefficients" = data.frame(Estimate, Std.Error),
      "deviance" = NULL,
      "vcov" = F.inv))
}

```

c)

Simulation of 1000 bernoulli draws with a random probability.

```
# probability
x = runif(1000, 0, 1)
# draw n bernoulli with prob x
y <- rbinom(1000, 1, x)
df <- data.frame(y, x)
### fit using glm
model <- glm(y ~ x, family = binomial(link = "probit"), data = df)
# beta
model$coefficients

## (Intercept)          x
##   -1.567633    3.135794

# se for beta
summary(model)

##
## Call:
## glm(formula = y ~ x, family = binomial(link = "probit"), data = df)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -2.2333  -0.7540   0.3502   0.7896   2.2322
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)  -1.5676     0.1002  -15.64  <2e-16 ***
## x              3.1358     0.1782   17.59  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 1386.3  on 999  degrees of freedom
## Residual deviance: 1002.6  on 998  degrees of freedom
## AIC: 1006.6
##
## Number of Fisher Scoring iterations: 4

# vcov
vcov(model)

##              (Intercept)          x
## (Intercept)  0.01004874 -0.01592284
## x            -0.01592284  0.03176926

# deviance
model$deviance

## [1] 1002.63

### fit using myglm
mymodel <- myglm(y ~ x, data = df)
# beta
```

```

mymodel$coefficients

##              Estimate Std. Error
## (Intercept) -1.567638  0.1002491
## x           3.135802  0.1782494

# vcov
mymodel$vcov

##              (Intercept)          x
## (Intercept)  0.01004989 -0.01592478
## x           -0.01592478  0.03177284

# deviance
mymodel$deviance

## NULL

```

Problem 2

a)

```

#install.packages("ISwR")
library(ISwR) # Install the package if needed
data(juul)
juul$menarche <- juul$menarche - 1
juul.girl <- subset(juul, age>8 & age<20 & complete.cases(menarche))

model <- glm(menarche ~ age, family=binomial(link="probit"), data= juul.girl)
anova(model, test = "Chisq")

## Analysis of Deviance Table
##
## Model: binomial, link: probit
##
## Response: menarche
##
## Terms added sequentially (first to last)
##
##
##      Df Deviance Resid. Df Resid. Dev  Pr(>Chi)
## NULL                    518      719.39
## age   1             522      197.39 < 2.2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

The low p-value suggests that age has an effect on the response variable.

b)

Relating to the juul data set, we define for each observation/individual

$$y_i = \begin{cases} 0, & \text{if menarche has occurred.} \\ 1, & \text{if menarche has not occurred.} \end{cases}$$

and t_i as the age at the time of examination, which corresponds to **age** in the data set. Let $T_i \sim N(\mu, \sigma)$, where T_i is the time until menarche occurs for the i 'th individual. Furthermore, let

$$\begin{aligned}\pi_i &:= P(y_i = 1) = P(T_i \leq t_i) \\ &= P\left(\frac{T_i - \mu}{\sigma} \leq \frac{t_i - \mu}{\sigma}\right) = \Phi\left(\frac{t_i - \mu}{\sigma}\right)\end{aligned}$$

This, in turn, gives

$$\Phi^{-1}(\pi_i) = -\frac{\mu}{\sigma} + \frac{1}{\sigma}t_i = \beta_0 + \beta_1 t_i,$$

where $\beta_0 = -\mu/\sigma$ and $\beta_1 = 1/\sigma$.