

Inconsistencies in the reproduction of figures 5.12 and 5.15 of Hosmer et al. 2013

Johann Popp

11 September 2017

Hosmer et al¹ have suggested several graphics to identify and investigate extreme and influencing covariate patterns in logistic regression models. These plots were reproduced using the computer software R², taking data from package `aplore3`³ and using package `epiR`⁴ for the calculation of statistics based on covariate patterns rather than single data rows.

```
par(mfrow = c(2,2))

# Load example data
glow <- aplore3::glow500

# Recode RATERISK
glow$raterisk3 <- cut(as.numeric(glow$raterisk), 2)
levels(glow$raterisk3) <- c("less/same", "greater")

# Logistic model from Table 4.16
model <- glm(fracture ~ age + height + priorfrac + momfrac + armassist + raterisk3 + age:priorfrac + mom:armassist, data=glow)

### Convert to covariate patterns
library(epiR)

# aggregate to covariate pattern
cp <- epi.cp(model.frame(model)[-1])
# Number of outcome events per covariate pattern
obs <- as.vector(by(as.numeric(as.factor(model$model[,1]))-1, as.factor(cp$id), sum))
# Estimated outcome probability per covariate pattern
fit <- as.vector(by(model$fitted.values, as.factor(cp$id), min))
# Calculate residuals and influence measures per covariate pattern
res <- epi.cpresids(obs, fit, cp)
# Calculate Change in deviance
deltaDeviance <- res$deviance^2 + ((res$pearson^2*res$leverage)/(1 - res$leverage))
# Put it together in one data.frame
dat <- data.frame(res, deltaDeviance, fit, cp$cov.pattern)

# Plot leverage vs. fit
plot(dat$fit, dat$leverage, pch = 16, cex = 0.8, xlab = "Estimated probability", ylab = "Leverage", xlim = c(0, 1), ylim = c(0, 1))
```

¹Hosmer, David W., Stanley Lemeshow, und Rodney X. Sturdivant. Applied logistic regression. 3rd Ed. Wiley series in probability and statistics. Hoboken, NJ: Wiley, 2013, p 186 ff.

²R Core Team (2017). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL <https://www.R-project.org/>.

³Luca Braglia (2016). `aplore3`: Datasets from Hosmer, Lemeshow and Sturdivant, “Applied Logistic Regression” (3rd Ed., 2013). R package version 0.9. <https://CRAN.R-project.org/package=aplore3>

⁴Mark Stevenson with contributions from Telmo Nunes, Cord Heuer, Jonathon Marshall, Javier Sanchez, Ron Thornton, Jeno Reiczigel, Jim Robison-Cox, Paola Sebastiani, Peter Solymos, Kazuki Yoshida, Geoff Jones, Sarah Pirikahu, Simon Firestone and Ryan Kyle. (2017). `epiR`: Tools for the Analysis of Epidemiological Data. R package version 0.9-87. <https://CRAN.R-project.org/package=epiR>

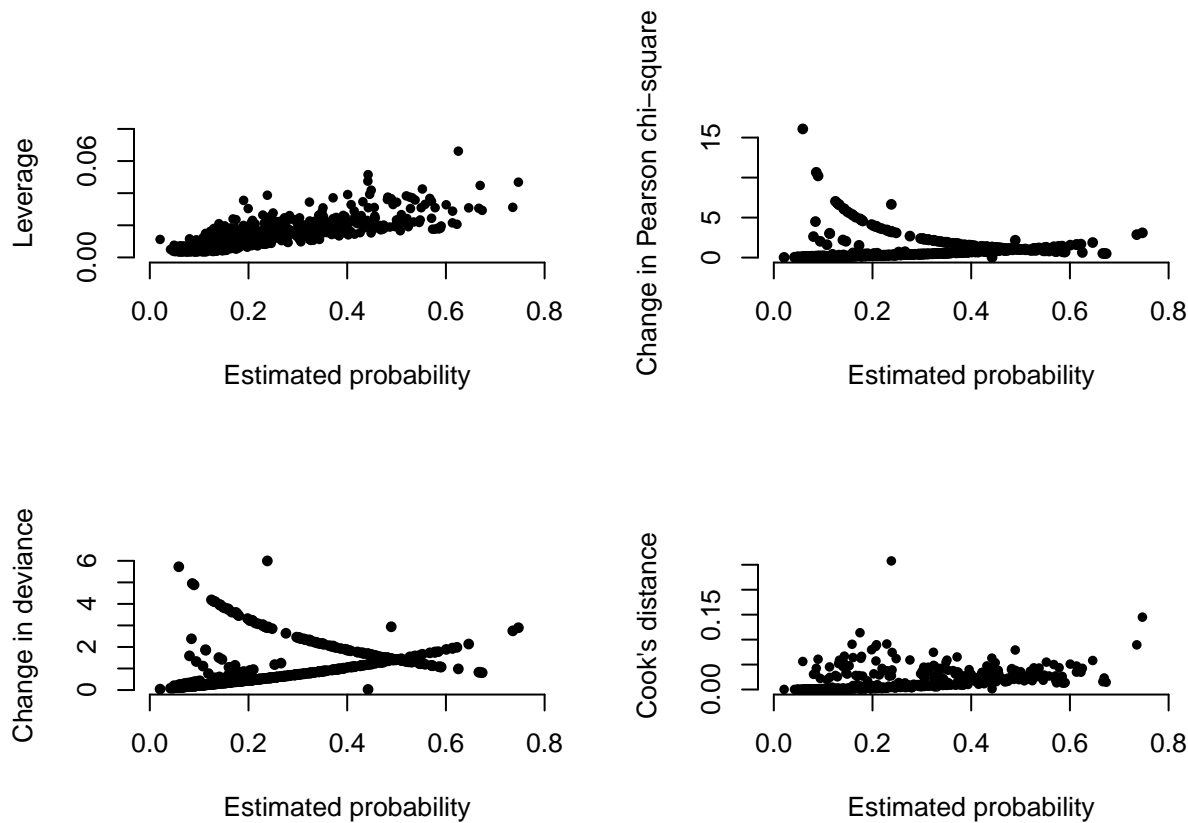
```

# Plot change in Chi-square vs. fit.
plot(dat$fit, dat$deltachi, pch = 16, cex = 0.9, xlab = "Estimated probability", ylab = "Change in Pearson chi-square")

# Plot change in deviance vs. fit.
plot(dat$fit, dat$deltaDeviance, pch = 16, cex = 0.9, xlab = "Estimated probability", ylab = "Change in deviance")

# Plot Cook's distance vs. fit
plot(dat$fit, dat$deltabeta, pch = 16, cex = 0.8, xlab = "Estimated probability", ylab = "Cook's distance")

```



```

par(mfrow = c(1,1))

```

This works fine for the graphs showing change of Pearson chi-square and change of deviance but the plots of leverage vs. fitted values and Cook's distance vs. fitted values differ substantially from figures 5.12 and 5.15 of the book.

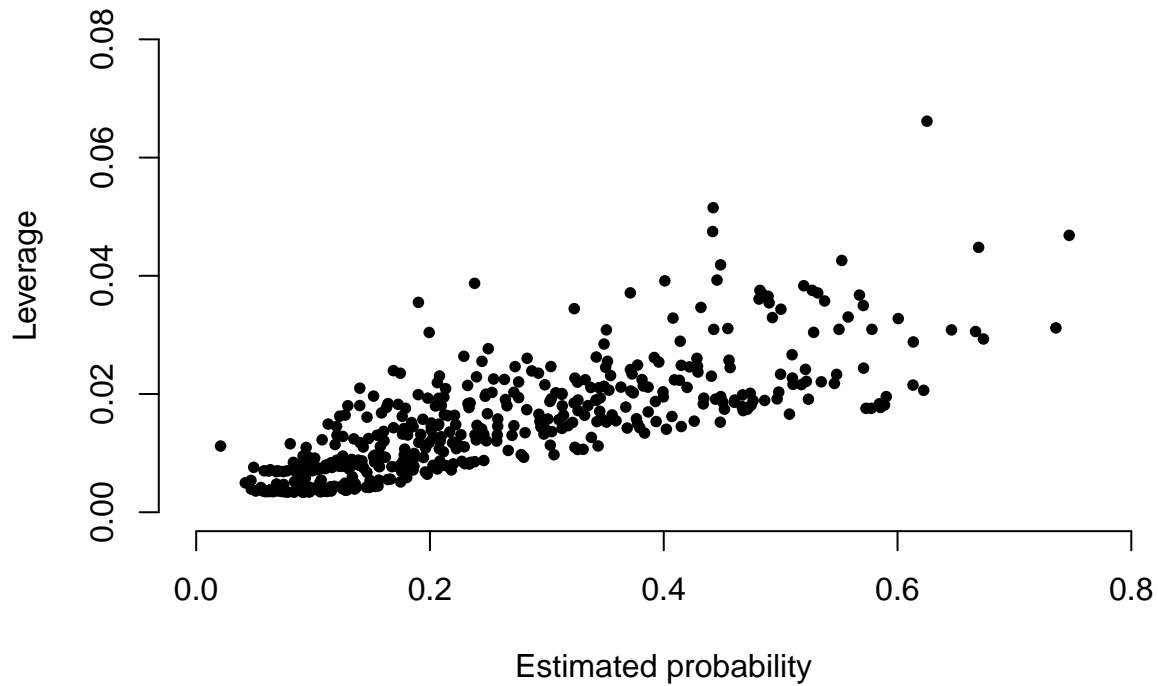
In contrast to what is said in the book, figure 5.12 of the book can be reproduced with a plot based on n-statistics (all the cases) instead of m-statistics (covariate patterns).

```

# Plot leverage vs. fit for covariate patterns
plot(dat$fit, dat$leverage, pch = 16, cex = 0.8, xlab = "Estimated probability", ylab = "Leverage", main = "Leverage vs. fit")

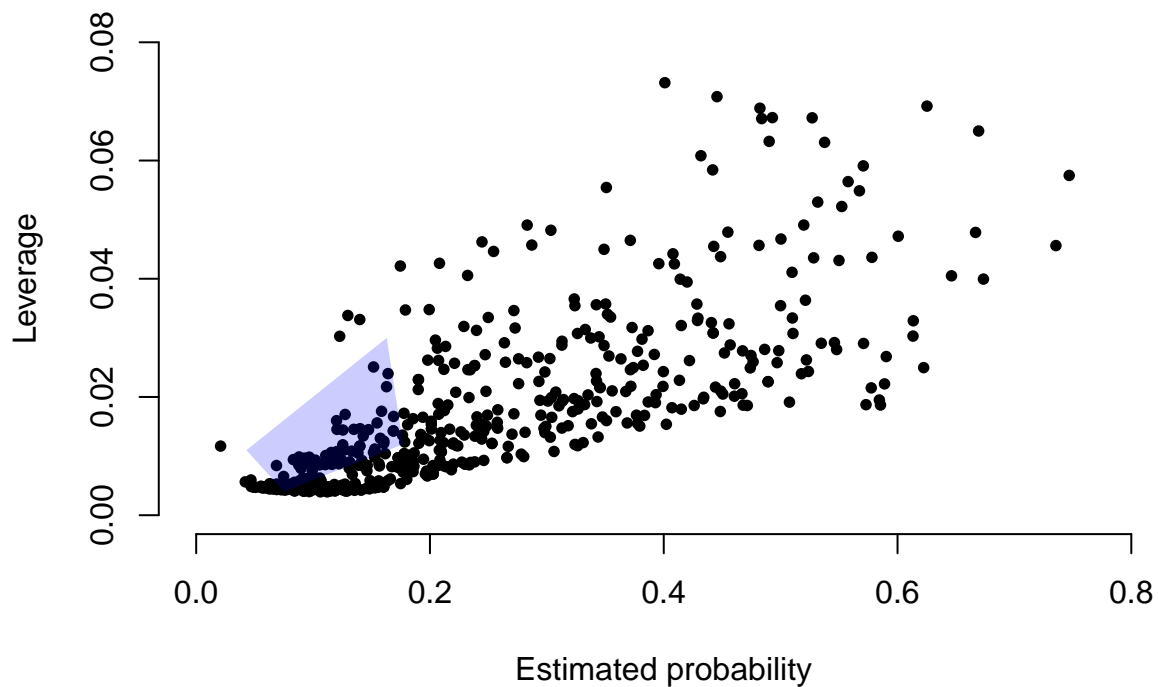
```

Leverage vs. fitted values based on covariate patterns



```
# Plot leverage vs. fit for each case
plot(fitted.values(model), hatvalues(model), pch = 16, cex = 0.8, xlab = "Estimated probability", ylab = "Leverage")
polygon(x = c(0.043, 0.075, 0.177, 0.163), y = c(0.011, 0.004, 0.012, 0.030), col = rgb(0,0,1,0.2), border = "black")
```

Leverage vs. fitted values based on individual cases



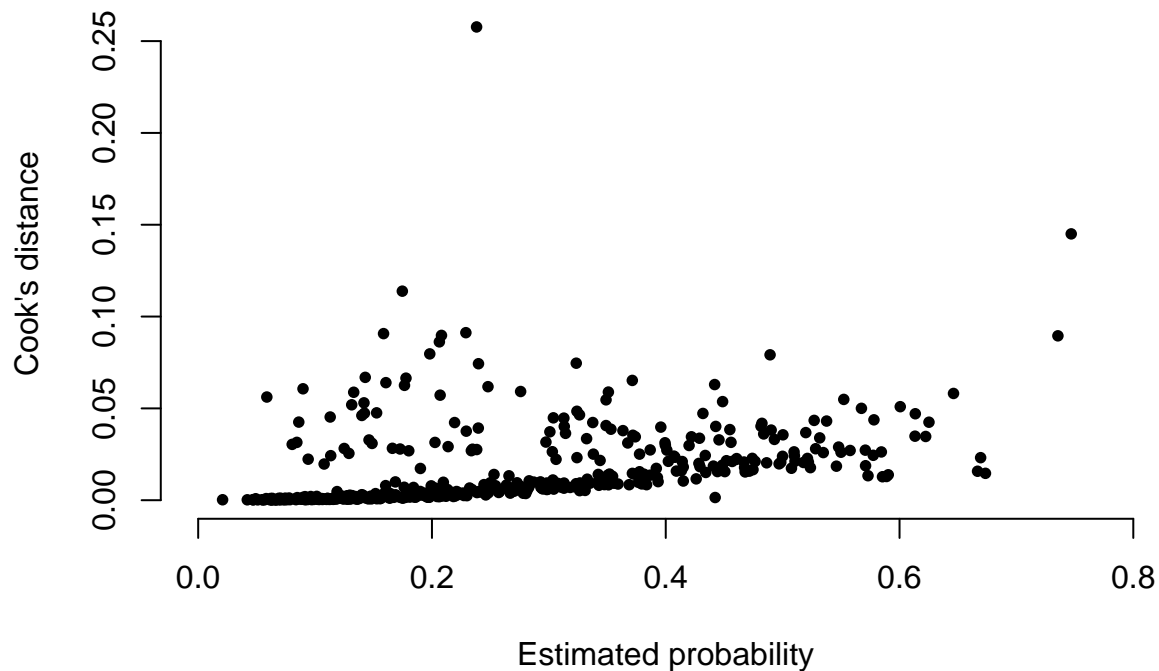
But still a close inspection of the highlighted area shows, that even this plot is not exactly the same as figure 5.12 in the book.

Regarding figure 5.15 some parts seem to be reproduced by a plot based on covariate patterns and other parts are quite well fitted by the reproduction based on individual cases. Overall the reproduced values of Cook's distance are substantially smaller than those shown in figure 5.15. Especially the highlighted point seems to be the one with a value of about 3.8 in the book.

```
# Plot Cook's distans vs. fit
```

```
plot(dat$fit, dat$deltabeta, pch = 16, cex = 0.8, xlab = "Estimated probability", ylab = "Cook's distans")
```

Cook's distance based on covariate patterns

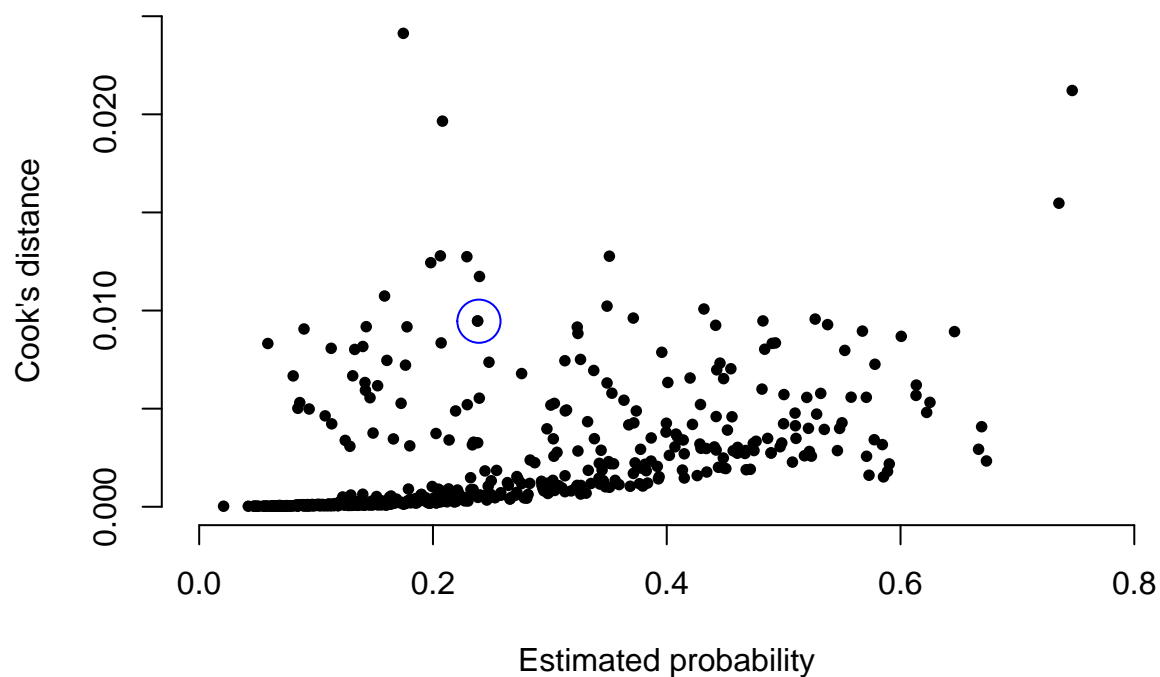


```
# Plot Cook's distans vs. fit
```

```
plot(fitted.values(model), cooks.distance(model), pch = 16, cex = 0.8, xlab = "Estimated probability", ylab = "Cook's distance")
```

```
points(x = 0.2393311, y = 0.009456836, col = "blue", cex = 3)
```

Cook's distance based on individual cases



As Cook's distance is based on leverage values it seems to me, that the calculation of those values in the package `epiR` differ from the calculation underlying the the figures of the book. As I am unfortunately not really into matrix algebra I can not effectively compare the syntax of the function `epiR::epi.cpresids` with formulas given in the book.