# R-bloggers

R news and tutorials contributed by (573) R bloggers

- Home
- About
- RSS
- add your blog!
- R jobs���
- Contact us

## Welcome!

Follow @rbloggers    25K f

Here you will find daily
**news and tutorials
about R**, contributed by
over 573 bloggers.
There are many ways to
**follow us** -
By e-mail:

Your e-mail here

Subscribe

21406 readers
BY FEEDBURNER

On Facebook:

R blog…
28k likes

Like Page

Be the first of your friends
to like this

**If you are an R blogger
yourself** you are invited
to add your own R
content feed to this site
(**Non-English** R bloggers
should add themselves-
here)

## Jobs for R-users

- Looking for a
  Freelance R-
  Developer with
  Shiny experience
- Data Analyst @
  New York
- R/Shiny App with
  d3 (small job, quick
  turnaround, $250 <
  4hrs)
- Postdoctoral
  position (Duke's
  River center) in
  data visualization
  and ecosystem
  science @ Durham,
  North Carolina,
  United States
- Software Engineer
  for The
  Computational
  Biology Program at
  Oregon Health and
  Science University
  @ Portland,
  Oregon, United
  States

Search & Hit Enter

# Popular Searches

- heatmap
- web scraping
- maps
- hadoop
- twitter
- alt=
- boxplot
- shiny
- time series
- animation
- ggplot2
- ggplot
- how to import image file to R
- latex
- trading
- finance
- PCA
- excel
- googlevis
- eclipse
- quantmod
- rstudio
- market research
- rattle
- Tutorial
- coplot
- rcmdr
- knitr
- title=
- rbloggers

# Recent Posts

- Partools, Recommender Systems and More
- Using htmlwidgets with knitr and Jekyll
- Wind in Netherlands
- Bioenergetics in R Workshop
- Correlation and Linear Regression
- Linear model with time series random component
- James Bond movies
- What it means to be a US Veteran Today
- Blog Post at Pluralsight
- The Lady Loves Statistics
- Annotables: R data package for annotating/converting Gene IDs
- Szkolenie z analizy sieciowej
- Applied Statistical Theory: Quantile Regression
- Let's meet on SatRdays: the link between RUGs and conferences
- importance sampling with infinite variance

# Other sites

- SAS blogs
- Statistics of Israel
- Jobs for R-users

# Veterinary Epidemiologic Research: GLM – Evaluating Logistic Regression Models (part 3)

March 19, 2013
By denishaine

Like    Share  〈 5    Tweet 〈 2

(This article was first published on **denis haine » R**, and kindly contributed to R-bloggers)

Third part on logistic regression (first here, second here).
Two steps in assessing the fit of the model: first is to determine if the model fits using summary measures of goodness of fit or by assessing the predictive ability of the model; second is to deterime if there's any observations that do not fit the model or that have an influence on the model.

**Covariate pattern**
A covariate pattern is a unique combination of values of predictor variables.

```
1   mod3 <- glm(casecont ~ dcpct + dneo + d
2   +              family = binomial("logit")
3   summary(mod3)
4
5   Call:
6   glm(formula = casecont ~ dcpct + dneo +
7       family = binomial("logit"), data = r
8
9   Deviance Residuals:
10      Min       1Q   Median       3Q
11   -1.9191  -0.7682   0.1874   0.5876   2.(
12
13   Coefficients:
14                  Estimate Std. Error z
15   (Intercept)   -3.776896   0.993251   -
16   dcpct          0.022618   0.007723
17   dneoYes        3.184002   0.837199
18   dcloxYes       0.445705   1.026026
19   dneoYes:dcloxYes -2.551997 1.205075   -
20   ---
21   Signif. codes:  0 '***' 0.001 '**' 0.01
22
23   (Dispersion parameter for binomial famil
24
25       Null deviance: 149.72  on 107  degre
26   Residual deviance: 103.42  on 103  degre
27   AIC: 113.42
28
29   Number of Fisher Scoring iterations: 5
30
31   library(epiR)
32   Package epiR 0.9-45 is loaded
33   Type help(epi.about) for summary informa
34
35   mod3.mf <- model.frame(mod3)
36   (mod3.cp <- epi.cp(mod3.mf[-1]))
37   $cov.pattern
38      id  n dcpct dneo dclox
39   1   1  7     0   No    No
40   2   2 38   100  Yes    No
41   3   3  1    25   No    No
42   4   4  1     1   No    No
43   5   5 11   100   No   Yes
44   8   6  1    25  Yes   Yes
45   10  7  1    14  Yes    No
46   12  8  4    75  Yes    No
47   13  9  1    90  Yes   Yes
48   14 10  1    30   No    No
49   15 11  3     5  Yes    No
50   17 12  9   100  Yes   Yes
51   22 13  2    20  Yes    No
52   23 14  8   100   No    No
53   25 15  2    50  Yes   Yes
54   26 16  1     7   No    No
55   27 17  4    50  Yes    No
56   28 18  1    50   No    No
```

```
57 │ 31  19  1    30   Yes    No
58 │ 34  20  1    99   No     No
59 │ 35  21  1    99   Yes    Yes
60 │ 40  22  1    80   Yes    Yes
61 │ 48  23  1     3   Yes    No
62 │ 59  24  1     1   Yes    No
63 │ 77  25  1    10   No     No
64 │ 84  26  1    83   No     Yes
65 │ 85  27  1    95   Yes    No
66 │ 88  28  1    99   Yes    No
67 │ 89  29  1    25   Yes    No
68 │ 105 30  1    40   Yes    No
69 │
70 │ $id
71 │    [1]  1  2  3  4  5  1  1  6  5  7  5
72 │   [26] 16 17 18  1  2 19  2 14 20 21 12
73 │   [51] 14 12 11  5 15  2  8  2 24  2  2
74 │   [76]  2 25  2 17  2  2  2  2 26 27 13
75 │  [101]  2  2  2  2 30  2  2  5
```

There are 30 covariate patterns in the dataset. The pattern dcpct=100, dneo=Yes, dclox=No appears 38 times.

**Pearson and deviance residuals**

Residuals represent the difference between the data and the model. The Pearson residuals are comparable to standardized residuals used for linear regression models. Deviance residuals represent the contribution of each observation to the overall deviance.

```
1 │ residuals(mod3) # deviance residuals
2 │ residuals(mod3, "pearson") # pearson resi
```

**Goodness-of-fit test**

All goodness-of-fit tests are based on the premise that the data will be divided into subsets and within each subset the predicted number of outcomes will be computed and compared to the observed number of outcomes. The Pearson $\chi^2$ and the deviance $\chi^2$ are based on dividing the data up into the natural covariate patterns. The Hosmer-Lemeshow test is based on a more arbitrary division of the data.

The Pearson $\chi^2$ is similar to the residual sum of squares used in linear models. It will be close in size to the deviance, but the model is fit to minimize the deviance and not the Pearson $\chi^2$. It is thus possible even if unlikely that the $\chi^2$ could increase as a predictor is added to the model.

```
1 │ sum(residuals(mod3, type = "pearson")^2)
2 │ [1] 123.9656
3 │ deviance(mod3)
4 │ [1] 103.4168
5 │ 1 - pchisq(deviance(mod3), df.residual(mo
6 │ [1] 0.4699251
```

The p-value is large indicating no evidence of lack of fit. However, when using the deviance statistic to assess the goodness-of-fit for a nonsaturated logistic model, the $\chi^2$ approximation for the likelihood ratio test is questionable. When the covariate pattern is almost as large as N, the deviance cannot be assumed to have a $\chi^2$ distribution.
Now the Hosmer-Lemeshow test, usually dividing by 10 the data:

```
 1 │ hosmerlem <- function (y, yhat, g = 10)
 2 │ +   cutyhat <- cut(yhat, breaks = quant:
 3 │ +              include.lowest = TRUI
 4 │ +   obs <- xtabs(cbind(1 - y, y) ~ cutyl
 5 │ +   expect <- xtabs(cbind(1 - yhat, yhat
 6 │ +   chisq <- sum((obs - expect)^2 / expe
 7 │ +   P <- 1 - pchisq(chisq, g - 2)
 8 │ +   c("X^2" = chisq, Df = g - 2, "P(>Ch:
 9 │ + }
10 │ hosmerlem(y = nocardia$casecont, yhat =
11 │ Erreur dans cut.default(yhat, breaks = (
12 │   'breaks' are not unique
```

The model used has many ties in its predicted probabilities (too few covariate values?) resulting in an error when running the Hosmer-Lemeshow test. Using fewer cut-points (g = 5 or 7) does not solve the problem. This is a typical example when not to use this test. A better goodness-of-fit test than Hosmer-Lemeshow and Pearson / deviance $\chi^2$ tests is the le Cessie – van Houwelingen – Copas – Hosmer unweighted sum of squares test for global goodness of fit (also here) implemented in the rms package (but you have to implement your model with the lrm function of this package):

```
1 │ mod3b <- lrm(casecont ~ dcpct + dneo + d
```

```
2  +                  method = "lrm.fit", model
3  +                  linear.predictors = TRUE,
4  residuals(mod3b, type = "gof")
5  Sum of squared errors    Expected value
6           16.4288056              16.82350
7                       Z
8           -1.4219860               0.15503
```

The p-value is 0.16 so there's no evidence the model is incorrect. Even better than these tests would be to check for linearity of the predictors.
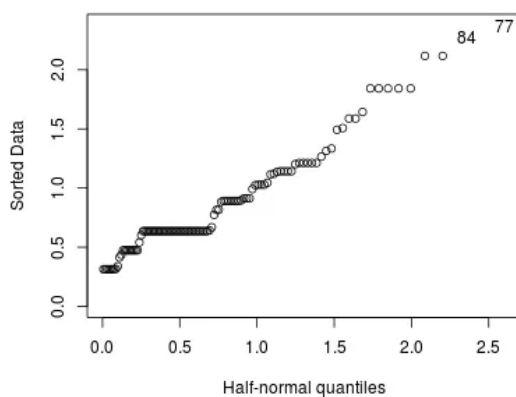
### Overdispersion

Sometimes we can get a deviance that is much larger than expected if the model was correct. It can be due to the presence of outliers, sparse data or clustering of data. The approach to deal with overdispersion is to add a dispersion parameter $\sigma^2$. It can be estimated with: $\hat{\sigma}^2 = \frac{X^2}{n-p}$ (p = probability of success). A half-normal plot of the residuals can help checking for outliers:

```
1  library(faraway)
2  halfnorm(residuals(mod1))
```



Half-normal plot of the residuals

The dispesion parameter of model 1 can be found as:

```
1   (sigma2 <- sum(residuals(mod1, type = "p
2   [1] 1.128778
3   drop1(mod1, scale = sigma2, test = "F")
4   Single term deletions
5
6   Model:
7   casecont ~ dcpct + dneo + dclox
8
9   scale:  1.128778
10
11          Df Deviance    AIC F value     Pr
12  <none>       107.99 115.99
13  dcpct    1   119.34 124.05 10.9350  0.001
14  dneo     1   125.86 129.82 17.2166 6.834e
15  dclox    1   114.73 119.96  6.4931  0.012
16  ---
17  Signif. codes:  0 '***' 0.001 '**' 0.01
18  Message d'avis :
19  In drop1.glm(mod1, scale = sigma2, test
20    le test F implique une famille 'quasil
```

The dispersion parameter is not very different than one (no dispersion). If dispersion was present, you could use it in the F-tests for the predictors, adding scale to drop1.
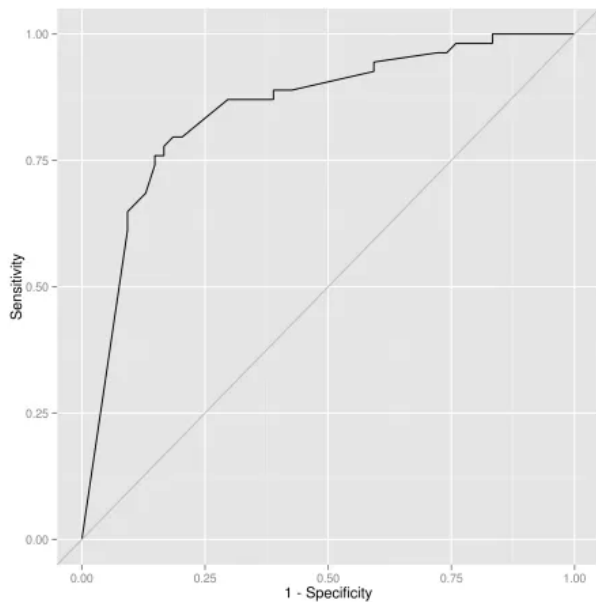
### Predictive ability of the model

A ROC curve can be drawn:

```
1   predicted <- predict(mod3)
2   library(ROCR)
3   prob <- prediction(predicted, nocardia$c
4   +                   label.ordering = c
5   tprfpr <- performance(prob, "tpr", "fpr"
6   tpr <- unlist(slot(tprfpr, "y.values"))
7   fpr <- unlist(slot(tprfpr, "x.values"))
8   roc <- data.frame(tpr, fpr)
9   ggplot(roc) + geom_line(aes(x = fpr, y =
10  +   geom_abline(intercept = 0, slope = 1
```

```
11  +      ylab("Sensitivity") +
12  +      xlab("1 - Specificity")
```
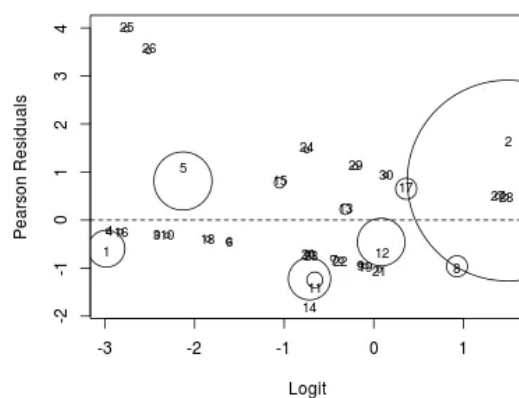


ROC curve

### Identifying important observations

Like for linear regression, large positive or negative standardized residuals
allow to identify points which are not well fit by the model. A plot of
Pearson residuals as a function of the logit for model 1 is drawn here, with
bubbles relative to size of the covariate pattern. The plot should be an
horizontal band with observations between -3 and +3. Covariate patterns
25 and 26 are problematic.

```
 1  nocardia$casecont.num <- as.numeric(noca
 2  mod1 <- glm(casecont.num ~ dcpct + dneo
 3  +            data = nocardia) # "logit'
 4  mod1.mf <- model.frame(mod1)
 5  mod1.cp <- epi.cp(mod1.mf[-1])
 6  nocardia.cp <- as.data.frame(cbind(cpid
 7  +                              no
 8  +                              fi
 9  ### Residuals and delta betas based on
10  mod1.obs <- as.vector(by(as.numeric(noc
11  +                       as.factor(no
12  mod1.fit <- as.vector(by(nocardia.cp$fi
13  +                       FUN = min))
14  mod1.res <- epi.cpresids(obs = mod1.obs
15  +                         covpattern =
16
17  mod1.lodds <- as.vector(by(predict(mod1
18  +                           FUN = min)
19
20  plot(mod1.lodds, mod1.res$spearson,
21  +       type = "n", ylab = "Pearson Resi
22  text(mod1.lodds, mod1.res$spearson, lab
23  symbols(mod1.lodds, mod1.res$spearson, c
```

Bubble plot of standardized residuals

The hat matrix is used to calculate leverage values and other diagnostic parameters. Leverage measures the potential impact of an observation. Points with high leverage have a potential impact. Covariate patterns 2, 14, 12 and 5 have the largest leverage values.

```
1  mod1.res[sort.list(mod1.res$leverage, de
2  cpid  leverage
3  2   0.74708052
4  14  0.54693851
5  12  0.54017700
6  5   0.42682684
7  11  0.21749664
8  1   0.19129427
9  ...
```

Delta-betas provides an overall estimate of the effect of the $j^{th}$ covariate pattern on the regression coefficients. It is analogous to Cook's distance in linear regression. Covariate pattern 2 has the largest delta-beta (and represents 38 observations).

```
1  mod1.res[sort.list(mod1.res$sdeltabeta,
2  cpid sdeltabeta
3  2    7.890878470
4  14   3.983840529
5  ...
```

Like    Share ⟨ 5     Tweet ⟨ 2

**Related**

$(75 - 50) \times 0.022 = 0.55$

[Veterinary Epidemiologic Research: GLM – Logistic Regression (part 2)](#)
In "R bloggers"

[Logistic Regression with R](#)
In "R bloggers"

[Veterinary Epidemiologic Research: GLM – Logistic Regression](#)
In "R bloggers"

5     2
Like     Tweet
Share

If you got this far, why not **subscribe for updates** from the site? Choose your flavor: e-mail, twitter, RSS, or facebook...

Like    Share ⟨ 5     Tweet ⟨ 2

Comments are closed.

# Top 3 Posts from the past 2 days

- [James Bond movies](#)
- [Correlation and Linear Regression](#)
- [Linear model with time series random component](#)

Search & Hit Enter

# Top 9 articles of the week

1. [Installing R packages](#)

## Sponsors

## ⬛ Jobs for R users

- Looking for a Freelance R-Developer with Shiny experience
- Data Analyst @ New York
- R/Shiny App with d3 (small job, quick turnaround, $250 < 4hrs)
- Postdoctoral position (Duke's River center) in data visualization and ecosystem science @ Durham, North Carolina, United States
- Software Engineer for The Computational Biology Program at Oregon Health and Science University @ Portland, Oregon, United States
- Data Scientist @ Israel
- Jr Financial Engineer

**Full list of contributing R-bloggers**

**R-bloggers** was founded by Tal Galili, with gratitude to the R community.
Is powered by WordPress using a bavotasan.com design.
Copyright © 2015 **R-bloggers**. All Rights Reserved. Terms and Conditions for this website