

Generalized linear models

Juan R Gonzalez

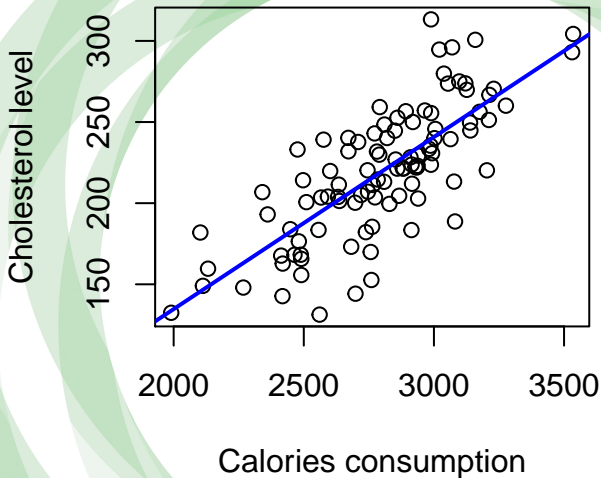
BRGE - Bioinformatics Research Group in Epidemiology
Barcelona Institute for Global Health (ISGlobal)
e-mail: juanr.gonzalez@isglobal.org
<http://brge.isglobal.org>
and Departament of Mathematics, UAB

- Linear regression
- Logistic regression
- Poisson and negative binomial regression
- Joinpoint and segmented regression
- Linear mixed models

Regression modeling

Outcome	Method	Example
Continuous	Linear regression	Factors that affects cholesterol levels
Binary	Logistic regression	Factors that affects developing cancer
Count	Poisson and Negative Binomial regression	Incidence and mortality trends
All	Joinpoint and segmented regression	Changes in longitudinal data
Time to event	Survival	Factors that affect time until developing cancer
All	Repeated/clustered measures	Factors that affect outcome complex data s

Linear regression



Linear regression

$$Y = \alpha + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n + \epsilon$$

- α correspond to the mean level of Y in the population
- β_j indicates the change in Y when X_j changes in 1 unit (after keeping the rest of X_k fixed)

Linear regression

Example: Researchers are interested in knowing the factors that better explain air Ozone levels (variable `Ozone` in data frame `airquality`). They measure solar radiation (`Solar.R`), average wind (`Wind`) and temperature (`Temp`) in different months (`Months`) for 154 observations.

```
data(airquality)
head(airquality)
```

##		Ozone	Solar.R	Wind	Temp	Month	Day
##	1	41	190	7.4	67	5	1
##	2	36	118	8.0	72	5	2
##	3	12	149	12.6	74	5	3
##	4	18	313	11.5	62	5	4
##	5	NA	NA	14.3	56	5	5
##	6	28	NA	14.9	66	5	6

Linear regression

Simple linear regression

```
mod <- lm(Ozone ~ Temp, data=airquality)
summary(mod)

##
## Call:
## lm(formula = Ozone ~ Temp, data = airquality)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -40.729 -17.409  -0.587   11.306  118.271
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -146.9955     18.2872  -8.038 9.37e-13 ***
## Temp          2.4287       0.2331  10.418 < 2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Linear regression

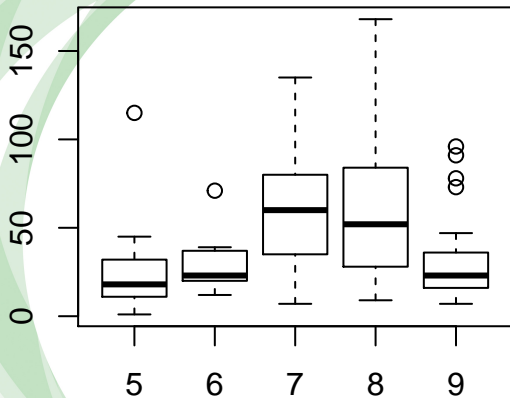
Multiple linear regression

```
mod <- lm(Ozone ~ Solar.R + Wind + Temp + as.factor(Month),
          data=airquality)
summary(mod)

##
## Call:
## lm(formula = Ozone ~ Solar.R + Wind + Temp + as.factor(Month),
##     data = airquality)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -40.344 -13.495  -3.165   10.399   92.689
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -74.23481    26.10184   -2.844   0.005
## Solar.R        0.05222     0.02367    2.206   0.029
## Wind         -3.10872     0.66009   -4.710 7.78e-
```


Linear regression

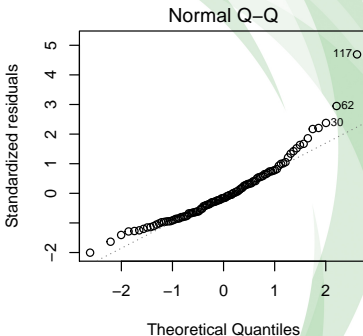
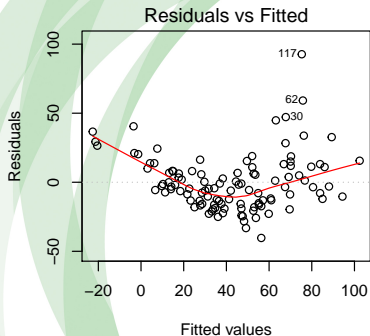
Interpretation of categorical factors



Linear regression

Model validation

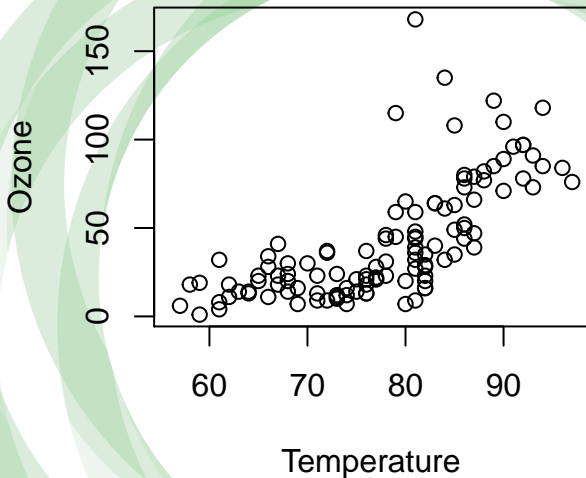
```
par(mfrow=c(2,2))  
plot(mod)
```



Scale-Location

Residuals vs Leverage

Linear regression



Linear regression

```
require(car)

## Loading required package: car
## Loading required package: carData

trans <- powerTransform(mod)
trans

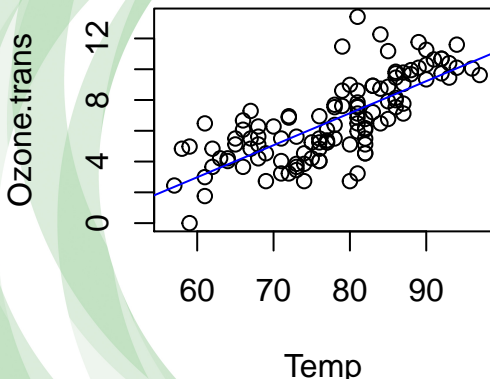
## Estimated transformation parameter
##          Y1
## 0.2206725

Ozone.trans <- bcPower(airquality$Ozone,
                      coef(trans, round=TRUE))

mod.trans <- lm(Ozone.trans ~ Temp, data=airquality)
```

Linear regression

```
plot(Ozone.trans ~ Temp, data=airquality)  
abline(mod.trans, col="blue")
```



Linear regression

Model validity can be measured by computing R^2

```
summary(mod)

##
## Call:
## lm(formula = Ozone ~ Temp, data = airquality)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -40.729 -17.409  -0.587   11.306  118.271
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -146.9955     18.2872  -8.038 9.37e-13 ***
## Temp          2.4287       0.2331   10.418 < 2e-16 ***
## ---
## Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
##
## Residual standard error: 23.71 on 114 degrees of freedom
## (37 observations deleted due to missingness)
## Multiple R-squared:  0.4877, Adjusted R-squared:  0.4832
## F-statistic: 108.5 on 1 and 114 DF,  p-value: < 2.2e-16
```

Linear regression

```
summary(mod.trans)
```

```
##
## Call:
## lm(formula = Ozone.trans ~ Temp, data = airquality)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.4144 -1.2733  0.0883  1.1028  6.0558
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -9.5085     1.3495  -7.046 1.49e-10 ***
## Temp           0.2082     0.0172  12.099 < 2e-16 ***
## ---
## Signif. codes:  0 *** 0.001 ** 0.01 * 0.05 . 0.1 1
##
## Residual standard error: 1.75 on 114 degrees of freedom
## (37 observations deleted due to missingness)
## Multiple R-squared:  0.5622, Adjusted R-squared:  0.5584
## F-statistic: 146.4 on 1 and 114 DF,  p-value: < 2.2e-16
```

Logistic regression

Y variable is binary (case/control, relapse/non-relapse, mortality, ...).
In that case, the logit transformation guarantees linearity.

$$\log(p(Y = 1)/(1 - p(Y = 1))) = \alpha + \beta_1 X_1 + \dots + \beta_k X_k$$

$\exp(\beta_k)$ can be interpreted as the odds ratio (OR) of
having/developing/being $Y = 1$

Logistic regression

Example: Reserchers are interested in determining whether a new treatment (variable `rx`) reduces mortality (variable `fustat`) in patients diagnosed with ovarian cancer. Data are available by typing:

```
data(ovarian, package="survival")
head(ovarian)
```

##	futime	fustat	age	resid.ds	rx	ecog.ps
## 1	59	1	72.3315	2	1	1
## 2	115	1	74.4932	2	1	1
## 3	156	1	66.4658	2	1	2
## 4	421	0	53.3644	2	2	1
## 5	431	1	50.3397	2	1	1
## 6	448	0	56.4301	1	1	2

Logistic regression

```
mod2 <- glm(fustat ~ rx, data=ovarian, family="binomial")
summary(mod2)

##
## Call:
## glm(formula = fustat ~ rx, family = "binomial", data = ovarian)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.2435  -0.9854  -0.9854   1.1127   1.3824
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)   0.7783     1.2502   0.623   0.534
## rx           -0.6242     0.7966  -0.784   0.433
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 35.890  on 25  degrees of freedom
## Residual deviance: 35.268  on 24  degrees of freedom
## AIC: 39.268
##
## Number of Fisher Scoring iterations: 4
```