

# Generalized linear models

Juan R Gonzalez

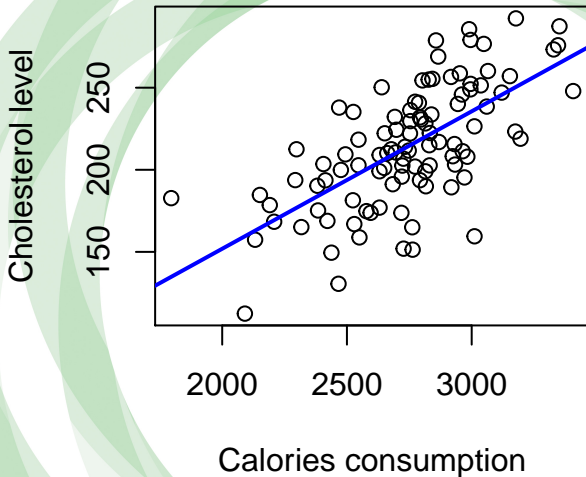
BRGE - Bioinformatics Research Group in Epidemiology  
Barcelona Institute for Global Health (ISGlobal)  
e-mail: [juanr.gonzalez@isglobal.org](mailto:juanr.gonzalez@isglobal.org)  
<http://brge.isglobal.org>  
and Departament of Mathematics, UAB

- Linear regression
- Logistic regression
- Poisson and negative binomial 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 sets

# Linear regression



# Linear regression

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

- $\alpha$  correspond to the mean level of  $Y$  in the population
- $\beta_j$  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

##

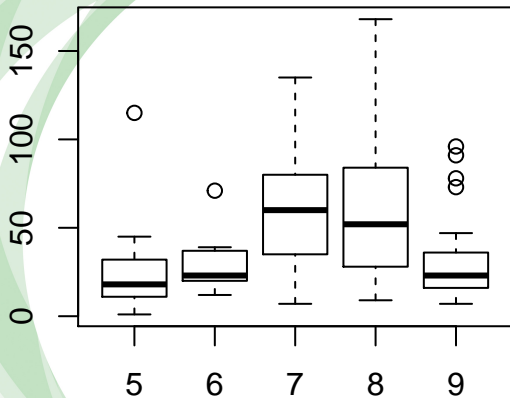
## Coefficients:

##	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

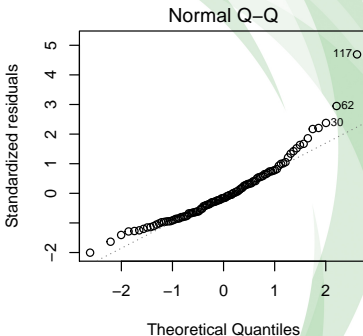
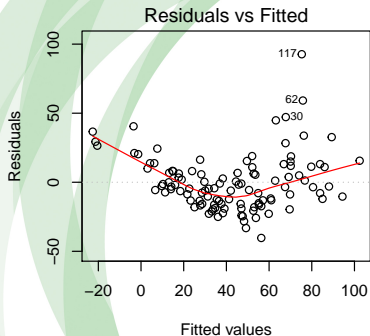
## Interpretation of categorical factors



# Linear regression

## Model validation

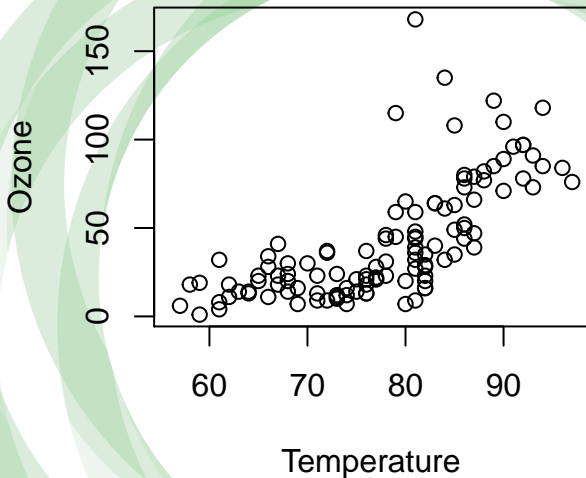
```
par(mfrow=c(2,2))  
plot(mod)
```



Scale-Location

Residuals vs Leverage

# Linear regression



# Linear regression

```
require(car)

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

trans <- powerTransform(mod)
trans

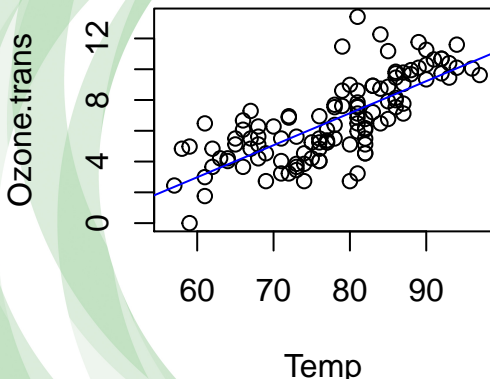
## Estimated transformation parameter
##          Y1
## 0.2206725

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
```

# Poisson regression

Let us analyze the following data that encodes Kentucky yearly cancer mortality from 1999-2005. Our aim is to study the evolution of mortality rates.

```
library(ljr)
```

```
## ljr 1.4-0 loaded
```

```
data(kcm)
```

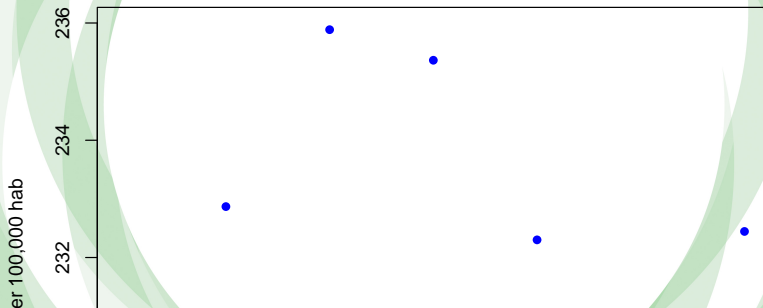
```
head(kcm)
```

##	Year	Count	Population
## 1	1999	9196	4018053
## 2	2000	9412	4041769
## 3	2001	9595	4067643
## 4	2002	9624	4088977
## 5	2003	9558	4114489
## 6	2004	9373	4140427

# Poisson regression

Let us visualize the evolution or mortality rates

```
kcm$ tasa <- (kcm$Count/kcm$Population)*100000  
plot(kcm$Year, kcm$tasa, xlab="Year",  
      ylab="Mortality rate per 100,000 hab",  
      type="n")  
points(kcm$Year, kcm$tasa, pch=16, col="blue")
```



# Poisson regression

This model evaluates the overall trend using a Poisson model:

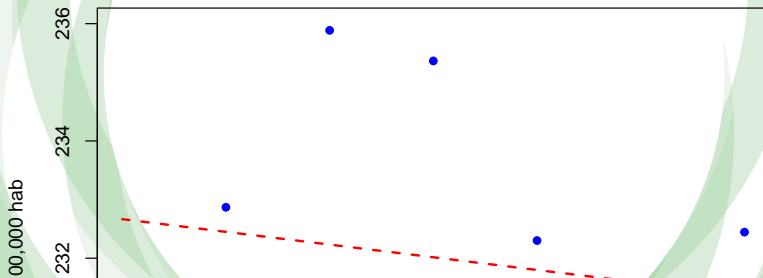
```
modPoisson <- glm(Count~Year+offset(log(Population)),  
                  family=poisson, data=kcm)  
modPoisson  
  
##  
## Call:  glm(formula = Count ~ Year + offset(log(Popula  
##      data = kcm)  
##  
## Coefficients:  
## (Intercept)          Year  
## -4.1972066    -0.0009335  
##  
## Degrees of Freedom: 6 Total (i.e. Null);  5 Residual  
## Null Deviance:      12.2  
## Residual Deviance: 11.96  AIC: 92.94
```

The percentage anual change is estimated by

```
round((1 - exp(modPoisson$coef[2]))*100, 2)
```

# Poisson regression

```
counts.pred <- predict(modPoisson, type="response")
tasa.pred <- (counts.pred/kcm$Population)*100000
plot(kcm$Year, kcm$tasa, xlab="Year",
     ylab="Mortality rate per 100,000 hab",
     type="n")
points(kcm$Year, kcm$tasa, pch=16, col="blue")
lines(kcm$Year, tasa.pred, lwd=2, lty=2, col="red")
```



# Poisson regression

Dispersion can be estimated using the residual deviance. A coefficient  $>1$  indicates that overdispersion is present. The a negative binomial regression is required. An approximated test can be (see others in the library `pscl`)

$H_0$  : There is no-overdispersion (1)

The associated p-value can be obtained by means of:

```
1 - pchisq(modPoisson$deviance, modPoisson$df.res)

## [1] 0.03527432

modPoisson

##
## Call: glm(formula = Count ~ Year + offset(log(Popula
##      data = kcm)
##
## Coefficients:
## (Intercept)          Year
```

# Negative Binomial regression

```
library(MASS)
modNB <- glm.nb(Count~Year+ offset(log(Population)),
               data=kcm)

modNB

##
## Call:  glm.nb(formula = Count ~ Year + offset(log(Population)),
##             init.theta = 13398.06067, link = log)
##
## Coefficients:
## (Intercept)          Year
## -4.2218596    -0.0009212
##
## Degrees of Freedom: 6 Total (i.e. Null);  5 Residual
## Null Deviance:      7.141
## Residual Deviance: 7.01  AIC: 93.73
```



# Regresión Joinpoint con R

En general, podemos estar interesados en estimar el mejor modelo para un número prefijado de joinpoint

Para 1 joinpoint

```
ljrjk(1, kcm$Count, kcm$Population, kcm$Year+.5)

## Model:
## y~Binom(n,p) where p=invlogit(eta)
## eta=b0+g0*t+g1*max(t-tau1,0)
##
##           Variables           Coef
## b0      Intercept -40.81272431
## g0              t    0.01737196
## g1 max(t-tau1,0)  -0.02418284
##
## Joinpoints:
##
## 1 tau1= 2001.273
## $Coef
##           Intercept           t max(t-tau1,0)
```