

# Generalized Linear Models

- Classical linear regression
  - ⇒ complicated formulation of simple model, structural and random component of the model
- Generalized linear models
  - ⇒ general description and examples
- Parameter estimation (extra material)<sup>1</sup>
  - ⇒ maximum likelihood method, computational issues
- Statistical inference (extra material)<sup>1</sup>
  - ⇒ goodness of fit, analysis of deviance

<sup>1</sup>Not obligatory

# Wikipedia

- In statistics, the **generalized linear model (GLM)** is a useful generalization of ordinary least squares regression. It relates the random distribution of the measured variable of the experiment (the *distribution function*) to the systematic (non-random) portion of the experiment (the *linear predictor*) through a function called the *link function*.
- The subject of generalized linear models was formulated by John Nelder and Robert Wedderburn as a way of unifying various other statistical models under one framework, **allowing for one general method of efficiently performing maximum likelihood estimation for these models.**

# Classical Linear Regression

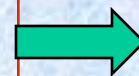
*Why easy formulation if complicated formulation exists?*

Response variable  $Y$  has a normal distribution

Expected value  $EY$  of  $Y$  depends on explanatory variables

- $(i) \quad Y_i \sim N(\mu_i, \sigma^2)$  *random component*
- $(ii) \quad \eta_i = x_i^T \beta$  *systematic component, linear predictor*
- $(iii) \quad \eta_i = g(\mu_i) = \mu_i$  *link function, linking (i) and (ii)*

Taking more general distribution in (i)  
and a more general link function  $g$  in  
(iii) sticking to this linear form in (ii)



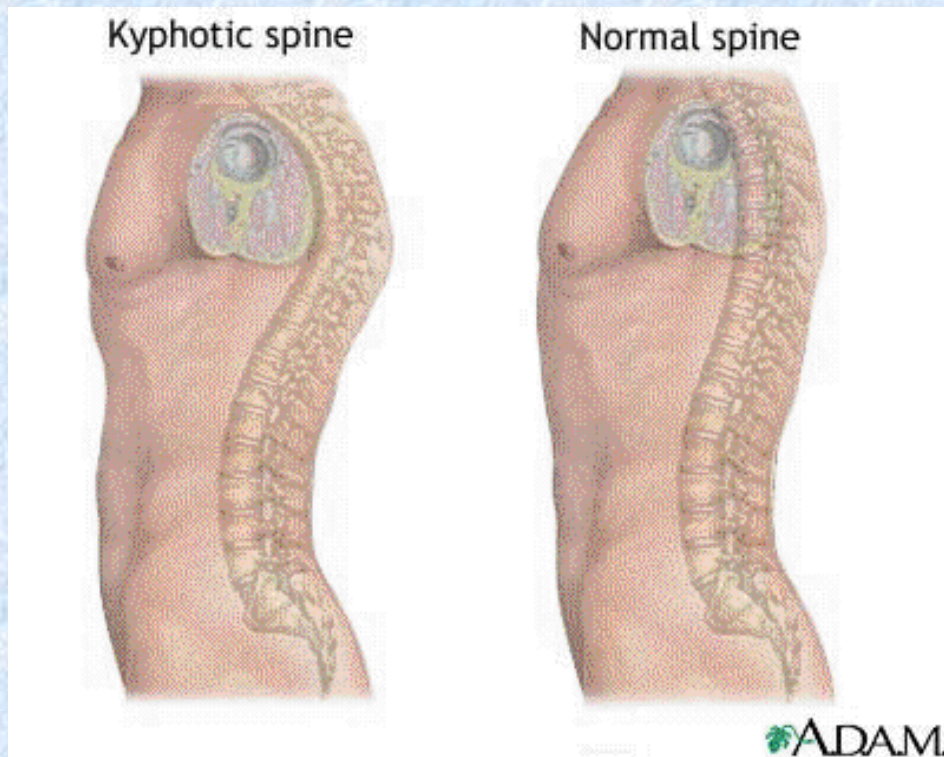
Generalized **L**inear **M**odel



# Generalized Linear Models

## *Kyphosis, medical context*

*Kyphosis* is a deformation that can occur with children that underwent corrective spinal surgery.



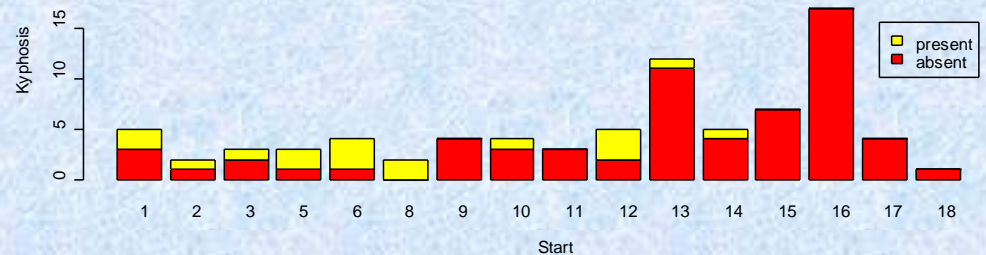
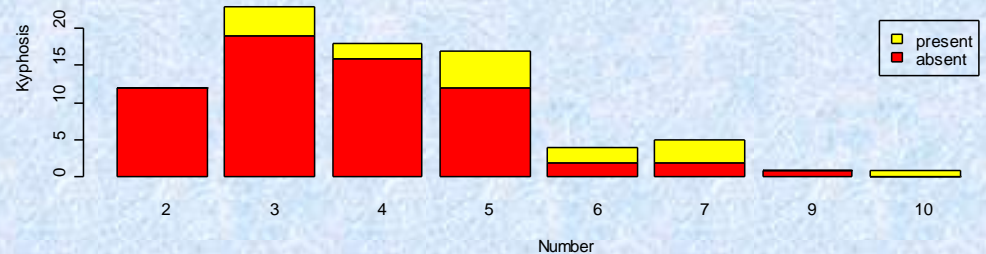
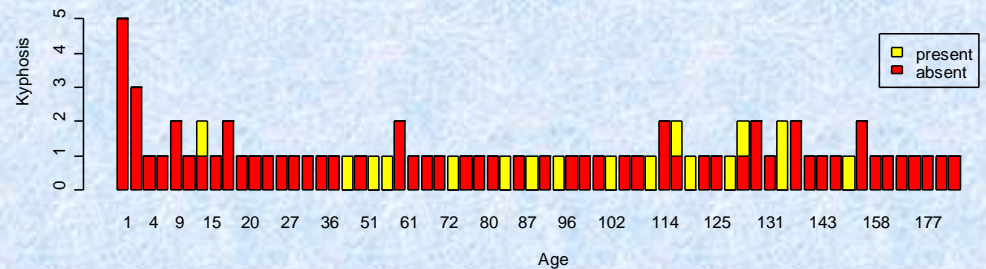
**Question:** given the age of the child at the time of surgery, what is the probability of occurrence of Kyphosis?

```
> library(rpart)
> data(kyphosis)
> kyphosis
  Kyphosis Age Number Start
1  absent  71      3      5
2  absent 158      3     14
...
```

# Kyphosis

## *Data visualisation*

```
> kyphosis
  Kyphosis Age Number Start
1  absent  71     3     5
2  absent 158     3    14
3  present 128     4     5
⋮      ⋮      ⋮      ⋮
83 absent  36     4    13
> par(mfrow=c(3,1))
> plot(Kyphosis~Age,
+ data=kyphosis)
> plot(Kyphosis~Number,
+ data=kyphosis)
> plot(Kyphosis~Start,
+ data=kyphosis)
```





# Kyphosis

*towards a model*

Notation:  $Y_i$  indicator of presence of Kyphosis (1 means present, 0 not)  
 $x_i$  age of child  $i$  at time of surgery

Initial naive model:

$$Y_i \sim \text{Bernoulli}(\mu_i), 1 \leq i \leq n, \text{ i.e., } Y_i = \begin{cases} 0 & \text{w.p. } 1 - \mu_i \\ 1 & \text{w.p. } \mu_i \end{cases}$$

$\mu_i$  depends on  $x_i$  in some way

$$Y_i \perp Y_j, i \neq j$$



# Kyphosis

*Generalized linear model*

Notation:  $Y_i$  indicator of presence of Kyphosis (1 means present, 0 not)  
 $x_i$  age of child  $i$  at time of surgery

→  $\mu_i$  does **not** depend on  $x_i$  **linearly**, but a nonlinear function  $g$  (**link function**) of  $\mu_i$  depends on  $x_i$  **linearly**

$$\left\{ \begin{array}{l} Y_i \overset{\text{indep}}{\sim} \text{Bernoulli}(\mu_i), 1 \leq i \leq n, \text{ i.e., } Y_i = \begin{cases} 0 & \text{w.p. } 1 - \mu_i \\ 1 & \text{w.p. } \mu_i \end{cases} \\ \quad \text{random component} \\ \eta_i = x_i^T \beta \quad \text{systematic component, linear predictor} \\ \eta_i = g(\mu_i) = \log(\mu_i / (1 - \mu_i)) \quad \text{link function (logit)} \end{array} \right.$$

$$\left( Y_i \overset{\text{indep}}{\sim} \text{Bern} \left( \frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)} \right), 1 \leq i \leq n \right) \quad \longrightarrow \quad \text{Logistic Regression model}$$

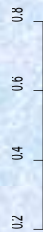
# Logistic Regression

$$\left\{ \begin{array}{l} Y_i \stackrel{\text{indep}}{\sim} \text{Bernoulli}(\mu_i), 1 \leq i \leq n, \text{ i.e., } Y_i = \begin{cases} 0 & \text{w.p. } 1 - \mu_i \\ 1 & \text{w.p. } \mu_i \end{cases} \\ \text{random component} \end{array} \right.$$

$$\eta_i = x_i^T \beta \quad \text{systematic component, linear predictor}$$

$$\eta_i = g(\mu_i) = \log(\mu_i / (1 - \mu_i)) \quad \text{link function}$$

Bernoulli  $\mu \uparrow$



$$Y_i \stackrel{\text{indep}}{\sim} \text{Bern}\left(\frac{\exp(x_i^T \beta)}{1 + \exp(x_i^T \beta)}\right), 1 \leq i \leq n$$

$x^T \beta \rightarrow$



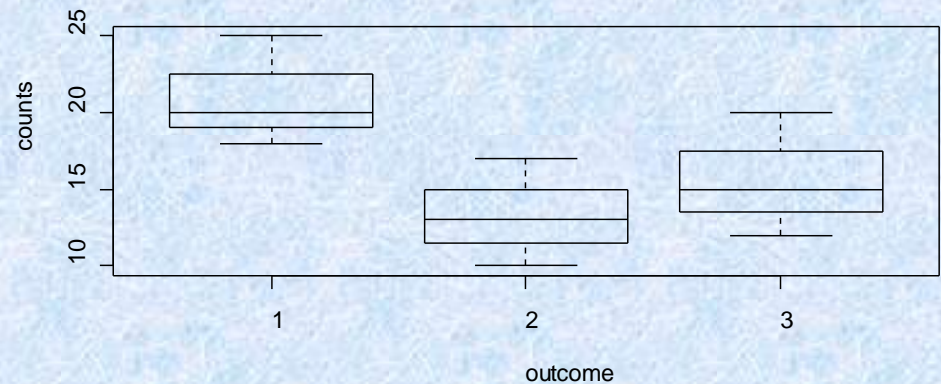
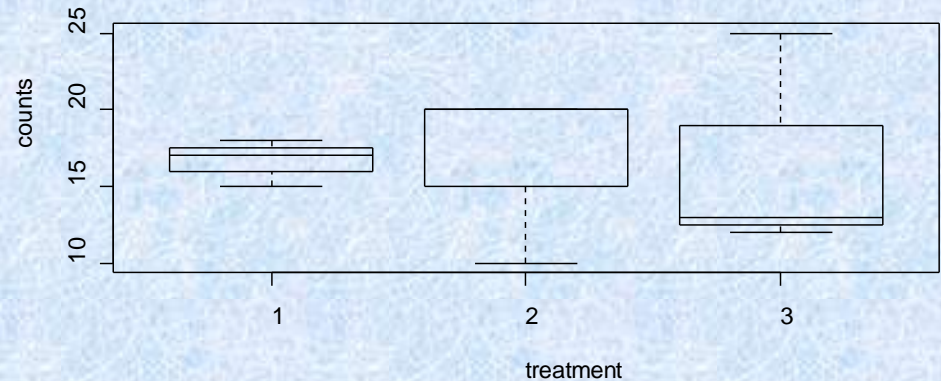
Estimate parameter vector  $\beta$  based on the available data  $\Rightarrow$  estimated model specified



# Generalized Linear Models

*next example*

```
> d.AD
  treatment outcome counts
1         1         1    18
2         1         2    17
3         1         3    15
4         2         1    20
5         2         2    10
6         2         3    20
7         3         1    25
8         3         2    13
9         3         3    12
> par(mfrow=c(2,1))
> plot(counts~treatment,d.AD)
> plot(counts~outcome,d.AD)
```



Data: randomized control trial.

Check first example for R function glm (type ?glm)

# Counts

*Generalized linear model*

Notation:  $Y_i$  measurement on counts, assumed to be **Poisson**  
 $x_i$  vector of explanatory variables for experiment  $i$

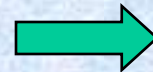
→  $\mu_i$  does **not** depend on  $x_i$  **linearly**, but a nonlinear function  $g$  (**link function**) of  $\mu_i$  depends on  $x_i$  **linearly**

$$\left\{ \begin{array}{ll} Y_i \overset{\text{indep}}{\sim} \text{Poisson}(\mu_i), 1 \leq i \leq n & \text{random component} \\ \eta_i = x_i^T \beta & \text{systematic component, linear predictor} \\ \eta_i = g(\mu_i) = \log \mu_i & \text{link function} \end{array} \right.$$

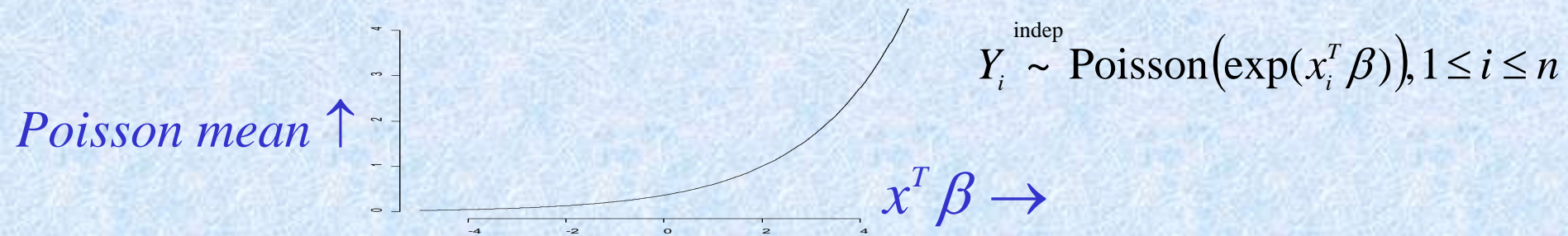
$$\left( Y_i \overset{\text{indep}}{\sim} \text{Poisson}(\exp(x_i^T \beta)), 1 \leq i \leq n \right) \longrightarrow \text{Log-linear Regression model}$$

# Log-linear Regression

$$\left\{ \begin{array}{l} Y_i \overset{\text{indep}}{\sim} \text{Poisson}(\mu_i), 1 \leq i \leq n \\ \eta_i = x_i^T \beta \\ \eta_i = g(\mu_i) = \log \mu_i \end{array} \right.$$



Estimate parameter vector  $\beta$  based on the available data  $\Rightarrow$  estimated model specified



How to estimate  $\beta$  in the logistic and log-linear regression model?



First: “general” generalized linear model  
 $\Rightarrow$  ML estimation in GLM's

# Generalized linear model

## *General structure and examples*

- (i)  $Y_i \sim f_i$ , a probability density function with  $EY_i = \mu_i$ ; see below
- (ii)  $\eta_i = x_i^T \beta$ , linear predictor
- (iii)  $\eta_i = g(\mu_i)$ , with  $g$  a general, monotonic link function

Here  $f_i$  is the probability density of a one-dimensional exponential family distribution

$$f_i(y) = \exp\left(\frac{y\theta_i - b(\theta_i)}{\phi/A_i} + c(y, \phi/A_i)\right)$$

$\theta_i$  natural parameter

$\phi$  dispersion parameter

Scale parameter, in statistical problems *known*

*Looks complicated, but.....*



# Generalized linear model

$$f_i(y) = \exp\left(\frac{y\theta_i - b(\theta_i)}{\phi/A_i} + c(y, \phi/A_i)\right)$$

*examples of exponential family*

$$\frac{1}{\sigma\sqrt{2\pi}} \exp(-(y - \mu)^2 / 2\sigma^2) = \exp\left(\frac{y\mu - \mu^2/2}{\sigma^2} - \frac{y^2}{2\sigma^2} - \log \sqrt{2\pi\sigma^2}\right)$$

$$\frac{\mu^y}{y!} \exp(-\mu) = \exp(y \log \mu - \log y! - \mu)$$
$$\binom{n}{y} \mu^y (1 - \mu)^{n-y} = \exp\left(y \log \frac{\mu}{1 - \mu} + n \log(1 - \mu) + \log \binom{n}{y}\right)$$

**Note:** in exponential family the **canonical link function** is the function mapping the mean  $\mu$  to the **natural parameter**. In other words: in a GLM with exponential family density  $f$  and canonical link, the natural parameter is modeled as linear function of the parameter vector  $\beta$ !!

# Generalized linear model

## Background material<sup>1</sup>

- Moments for exponential family
- Maximum Likelihood Estimation
- Newton-Rhapson
- Testing
- Confidence intervals

<sup>1</sup>Not obligatory