

Statistical Data Analysis

Dr. Jana de Wiljes

15. Dezember 2021

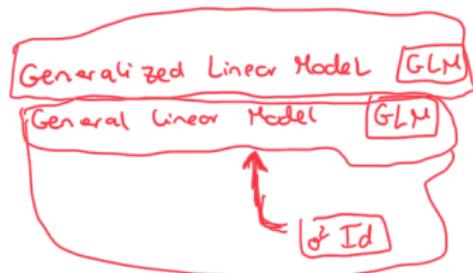
Universität Potsdam

Generalized Linear Models

Setting: $y_i = x_i^\top \beta + \epsilon, \quad i = 1, \dots, n$

Up till now:

- $\epsilon_i \sim N(0, \sigma^2)$
- $y_i \sim N(\mu_i, \sigma^2)$
- $\mu_i = x_i^\top \beta, \quad i = 1, \dots, n$
- $\mu_i = \mathbb{E}[y_i | x_i]$



Assumption:

$$f(y_i | x_i, \theta_i, \phi, w_i) = \exp\left(\frac{y_i \theta_i - b(\theta_i)}{a(\phi)} + c(y_i, \phi, w_i)\right) \quad (1)$$

where

- θ_i is the natural parameter of the family,
- ϕ is scale or dispersion parameter,
- $b(\cdot), c(\cdot), a(\cdot)$ are specific functions corresponding to the type of the family

$\phi = \sigma^2$ $a = id$

$M = x_i^\top \beta$

$\mu = h(x_i^\top \beta)$

How does the normal distribution fit in the picture?

$$\begin{aligned} f(y_i | x_i, \theta_i, \phi) &= \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y_i - \mu)^2}{2\sigma^2}\right) \\ &= \exp\left(\frac{y_i\mu - \mu^2/2}{\sigma^2} - \frac{y_i^2/\sigma^2 + \log(2\pi\sigma^2)}{2}\right) \end{aligned}$$

i.e.,

- $\theta_i = \mu$
- $a(\phi) = \sigma^2$
- $b(\theta_i) = \theta_i^2/2$
- $c(\cdot) = -\frac{1}{2}(y_i^2/\sigma^2 + \log(2\pi\sigma^2))$

$$[a - 2ab + b^2]$$

Mean and Variance?

Idea: • $\eta_i = g(\mu_i)$

↑

Link function (invertible) and nonlinear

• $\eta_i = \mathbf{x}_i^T \boldsymbol{\beta}$

↑

$\mu = g^{-1}(\eta_i) = g^{-1}(\mathbf{x}_i^T \boldsymbol{\beta})$

$g(\mu_i)$

↑

Mean and Variance?

Link mean to linear predictor

Naughty or nice ?



①
 $\gamma_i \approx 0.5$
②
0

Elves need system to
determine if child
naughty or nice

Idea:

- Blood measurements of
 1. Serontonixi
 2. Oximontiuos } x_i
- observation of a test group of kids for a year to identify label:
 1. naughty
 2. nice

Logistic regression

$$y_i \sim \text{Bin}(1, p_i) \quad i=1, \dots, n$$

$$y_i = \begin{cases} 1 & p_i \\ 0 & 1-p_i \end{cases}$$

$$f(y | \beta) = \prod_{i=1}^n p_i^{y_i} (1-p_i)^{1-y_i} = \prod_{i=1}^n \left(\frac{p_i}{1-p_i} \right)^{y_i} (1-p_i)$$

$$= \prod_{i=1}^n \exp \left(y_i \log \left(\frac{p_i}{1-p_i} \right) - \log \left(\frac{1}{1-p_i} \right) \right)$$

$$= \exp \left[\sum_{i=1}^n \{ y_i \theta_i - \log(1 + e^{\theta_i}) \} \right]$$

Logistic regression

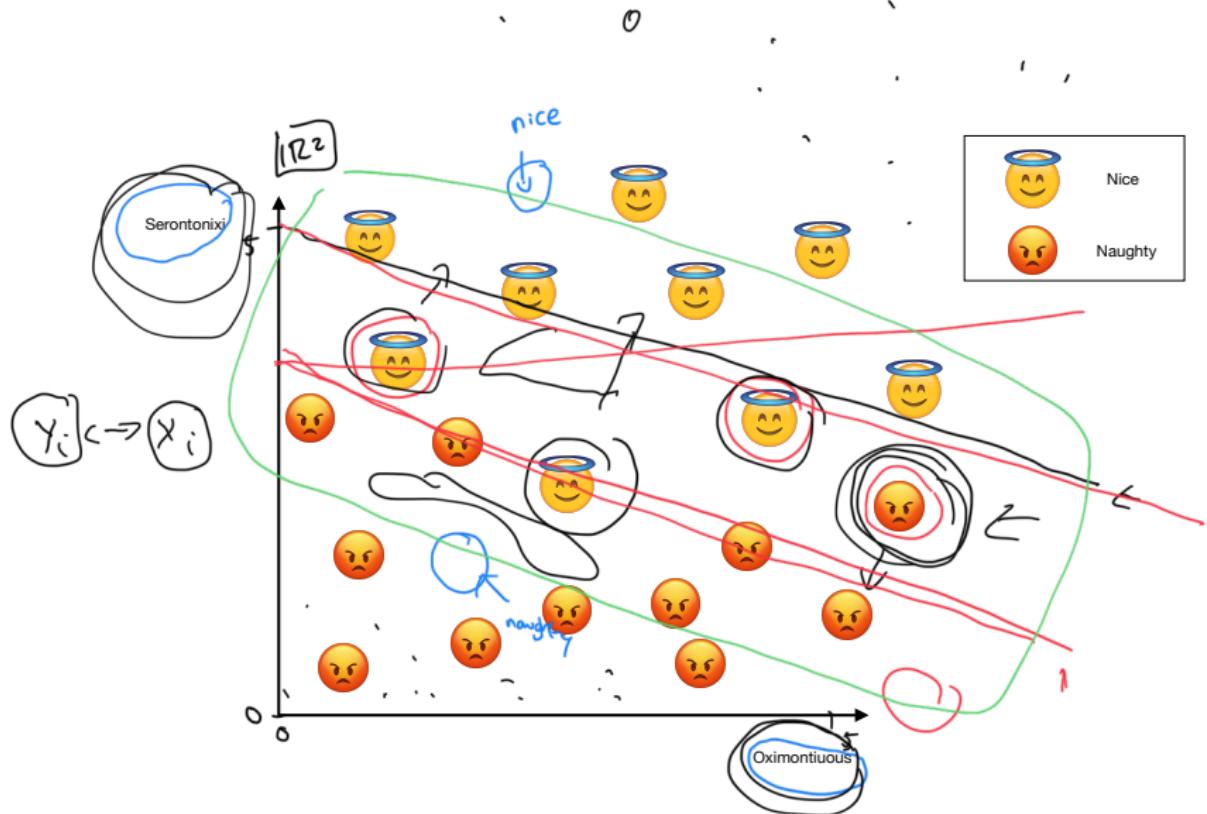
Different variations

	Normal	Poisson	Binomial	Gamma
Notation	$N(\mu_i, \sigma^2)$	$\text{Pois}(\mu_i)$	$\text{Bin}(n_i, \pi_i)$	$G(\mu_i, \nu)$
Range of y_i	$(-\infty, \infty)$	$[0, \infty)$	$[0, n_i]$	$(0, \infty)$
Dispersion, ϕ	σ^2	1	$1/n_i$	ν^{-1}
Cumulant: $b(\theta_i)$	$\theta_i^2/2$	$\exp(\theta_i)$	$\log(1 + e^{\theta_i})$	$-\log(-\theta_i)$
Mean function, $\mu(\theta_i)$	θ_i	$\exp(\theta_i)$	$1/(1 + e^{-\theta_i})$	$-1/\theta_i$
Canonical link: $\theta(\mu_i)$	identify	log	logit	reciprocal
Variance function, $V(\mu_i)$	1	μ	$\mu(1 - \mu)$	μ^2

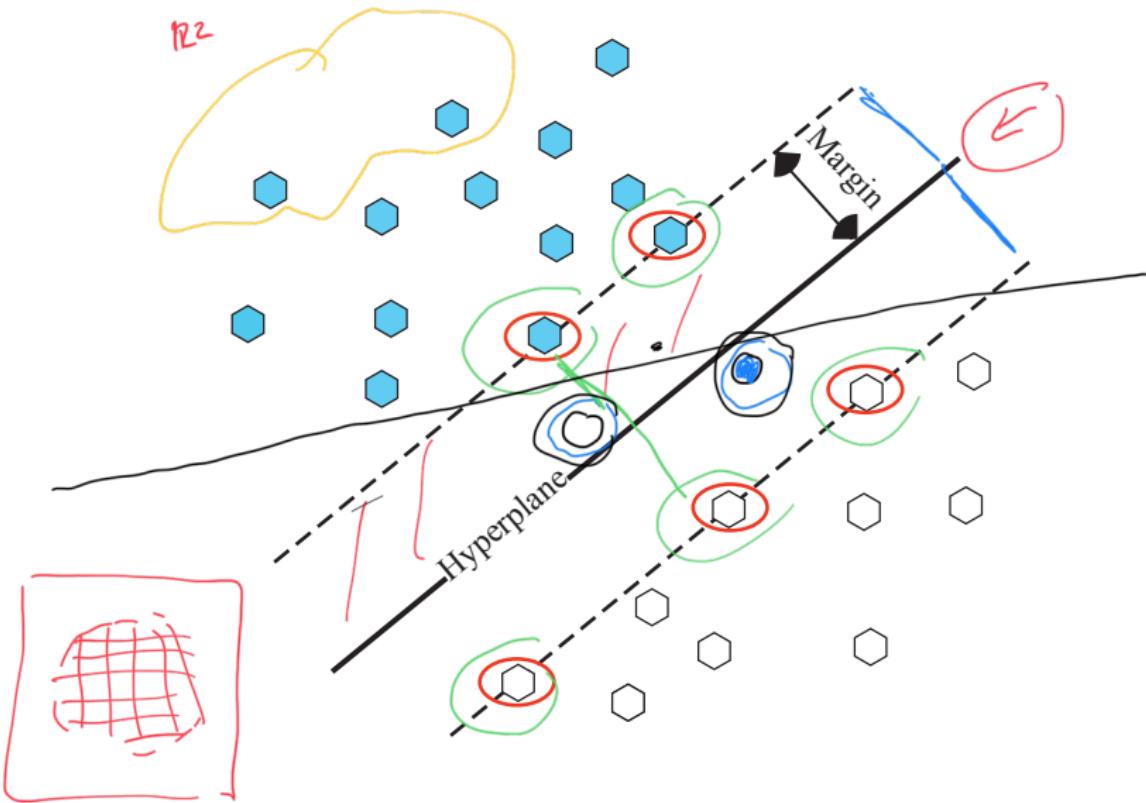


$$\text{logit } \eta_i = \log \left\{ \frac{\mu_i}{1 - \mu_i} \right\} = g(\mu_i)$$

Naughty or nice?



Support vector machine



Support vector machine

