

SEMESTER

2

19/6/16 ALM 2

Linear Regression

Linear regression and how it is generalised

Set of observations $\{(y_i, x_i)\}_{i=1, \dots, n}$ such that the below relationship holds

$$y_i = \beta_0 + \beta_1 x_{i1} + \dots + \beta_n x_{in} + \varepsilon_i \\ = \beta^T x_i + \varepsilon_i$$

Where $y_i \in \mathbb{R}$ outcome/response variable

$x_i = (1, x_{i1}, \dots, x_{in}) \in \mathbb{R}^{k+1}$ a vector collating values of the explanatory variable associated with the outcome y_i

ε_i is the noise, residual or error associated with y_i

$$\varepsilon_i \sim \mathcal{N}(0, \sigma^2) \quad \forall i \neq j \rightarrow \text{residual independent} \quad P(\varepsilon) = \frac{\exp(-\frac{\varepsilon^2}{2\sigma^2})}{\sqrt{2\pi}\sigma}$$

$P(y_i | x_i) \Rightarrow$ pdf or distribution

$$\mathbb{E}[y_i] = \mathbb{E}[x_i^T \beta + \varepsilon_i] = x_i^T \beta + 0$$

$$\text{Var}[y_i] = \sigma^2$$

$$P(y_i | \beta, x_i) = \mathcal{N}(y_i, \beta^T x_i, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}\right]$$

Best parameter β are estimated that maximise the joint probability density function of all the residuals: $\hat{\beta} = \arg\max_{\beta} p(\varepsilon_1, \dots, \varepsilon_n)$

Because residuals are independent and follow same distribution P_{ε} , the joint density function of residuals corresponds to

$$p(\varepsilon_1, \dots, \varepsilon_n) = \prod_{i=1}^n P_{\varepsilon}(\varepsilon_i) = \prod_{i=1}^n \frac{\exp(-\frac{\varepsilon_i^2}{2\sigma^2})}{\sqrt{2\pi}\sigma} = \prod_{i=1}^n \frac{\exp\left[-\frac{(y_i - \beta^T x_i)^2}{2\sigma^2}\right]}{\sqrt{2\pi}\sigma}$$

On redefining linear regression \hat{y}_i that:

- y_i is normally distributed

- With mean $\mathbb{E}[y_i] = \beta^T x_i$ and Variance $\mathbb{E}[(y_i - \beta^T x_i)^2] = \sigma^2$

19/01

Generalised Linear Models (GLM)

- Will generate three premises used for linear regression as follows:

• The probability density function $P(y|x, \beta)$ is a member of the exponential family of distributions. Other distributions available to deal with outcome y is not an element of \mathbb{R} but for instance is binary.

• Expectation of y given x and β is defined by:
 $E[y] = \int_{\mathcal{Y}} y P_{y|x, \beta}(y|x, \beta) dy$ with \mathcal{Y} domain of definition for outcome y

Now related to the exponential variational with a link function g such that $g[E[y]] = \beta^T x$

For instance, in the case of linear regression, the link function g that we used is the identity function g defined for $z \in \mathbb{R} \Rightarrow g(z) = z \in \mathbb{R}$

g is a link function that is bijective and its inverse g^{-1} exists.
 In general, this function maps the space of the expectation $E[y]$ to the space \mathbb{R} where $\beta^T x$ takes its value.

Probability Distributions

The pdf p_y of random variable y has the following properties:

• If $y \in \mathcal{Y}$ $p_y(y) \geq 0$: function p_y is positive for all possible outcomes
 \mathcal{Y} used as space of all possible outcomes

• The function p_y integrates to 1 on the space of all possible outcomes such that:
 - when y is a continuous r.v. $\int_{\mathcal{Y}} p_y(y) dy = 1$

- When y is discrete r.v. $\sum_{y \in \mathcal{Y}} p_y(y) = 1$

2 19/01/16 ALSM2

Exponential Family of Distribution

The distribution belongs to the exponential family if it can be written as:

$$p_{y|\theta}(y|\theta) = \exp[a(y)b(\theta) + c(\theta) + d(y)]$$

where a, b, c, d are known functions. If $a(y) = y$, then the distribution is said to be in canonical form.

Expectation

- y continuous $E[y] = \int_{\mathbb{R}} y p_y(y) dy$

- y discrete $E[y] = \sum_{y \in \mathbb{N}} y p_y(y)$

Normal / Gaussian Distribution

$$y \in \mathbb{R} \quad p_{y|\theta}(y|\theta, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{(y-\theta)^2}{2\sigma^2}\right] \quad \begin{matrix} \theta \in \mathbb{R} \\ \sigma^2 \in \mathbb{R}^+ \end{matrix} \quad E[y] = \theta$$

$$\exp\left[\frac{-(y-\theta)^2}{2\sigma^2} - \log(\sqrt{2\pi}\sigma)\right]$$

$$\exp\left[\underbrace{\frac{-y^2}{2\sigma^2}}_{d(y)} + \underbrace{\frac{2y\theta}{2\sigma^2}}_{a(y)b(\theta)} + \underbrace{\frac{\theta^2}{2\sigma^2} - \log(\sqrt{2\pi}\sigma)}_{c(\theta)}\right]$$

$$d(y) \quad a(y)b(\theta) \quad c(\theta)$$

Poisson Distribution

Expresses the probability of a given number of events y occurring in a fixed interval of time and/or space (and/or fixed total population size)

$$p_{y|\theta}(y|\theta) = \frac{\theta^y \exp(-\theta)}{y!} \quad y \in \mathbb{N} \quad \theta \in \mathbb{R}^+$$

$$E[y] = \sum_{y \in \mathbb{N}} y \frac{\theta^y}{y!} \exp(-\theta) = \theta$$

$$\exp(-\theta) \theta \sum_{y=1}^{\infty} \frac{\theta^{y-1}}{(y-1)!} \quad \text{Taylor expansion}$$

$$= \exp(-\theta) \exp(\theta) = \theta$$

θ is the parameter used to plug in information from explanatory variable x

$\theta \Rightarrow (x^T \beta)$ Somehow associated

$$= \frac{1}{x!} \exp[xy\lambda - \lambda]$$

$$= \frac{1}{x!} \exp(xy\lambda) \exp(-\lambda)$$

\uparrow \uparrow \uparrow
 $d(y)$ $a(y)b(\theta)$ $c(\theta)$

Bernoulli

y is binary variable, $P_{y|\theta}(y=1|\theta) = \theta$ $P_{y|\theta}(y=0|\theta) = 1-\theta$
 $P_{y|\theta} = \theta^y(1-\theta)^{1-y}$ $y \in \{0,1\}$ $\theta \in [0,1]$ $E[y] = \theta$

Binomial

Response y is # success in n trials, proportion θ real number $0,1$

$$P_{y|\theta}(y|\theta) = \frac{n!}{(n-y)!y!} \theta^y(1-\theta)^{n-y} \Rightarrow \binom{n}{y} \theta^y(1-\theta)^{n-y} \quad E[y] = n\theta$$

Exponential

$$P_{y|\theta}(y|\theta) = \theta \exp(-\theta y) \quad y \in \mathbb{R}^+ \quad \theta \in \mathbb{R}^{+\ast} \quad \leftarrow \text{and not } 0 \quad E[y] = \frac{1}{\theta}$$

$\exp(-\theta y + \log \theta)$
 \uparrow \uparrow \uparrow
 $b(y)$ $a(\theta)$ $d(y)=0$

$$\int_0^{+\infty} P(y|\theta) dy = 1 \quad \int_0^{+\infty} \theta \exp(-\theta y) dy : \left[-\exp^{-\theta y} \right]_0^{+\infty}$$

$$-\exp(-\infty) + \exp(0) = 0+1 = 1$$

Weibull

$$P_{y|\lambda, \theta}(y|\lambda, \theta) = \lambda \theta y^{\lambda-1} \exp[-\theta y^\lambda] \quad y \in \mathbb{R}^+ \quad \theta \in \mathbb{R}^{+\ast} \quad \lambda \in \mathbb{R}^{+\ast}$$

Expectation is $E[y] = \left(\frac{1}{\theta}\right)^{\frac{1}{\lambda}} \sim \left(1 + \frac{1}{\lambda}\right)$

$$\Gamma(u) = \int_0^{+\infty} s^{u-1} \exp(-s) ds$$

Exp is a special case of Weibull with $\lambda=1$

2 19/01/16 ALM 2

GLM

1. We have collected independently a set of responses y_i as well as the values for some explanatory variables stored in vector x_i i.e. have observations $\{(y_i, x_i)\}_{i=1, \dots, n}$
2. Response y_i has a distribution $p_{y_i}(\theta_i | \theta_i)$ that is a member of exponential family, indexed by parameter θ and related to expectation of response $E[y_i]$
3. Model constructed by linking expectation of response $E[y_i]$ with the linear predictor $x_i^T \beta$
$$g(E[y_i]) = x_i^T \beta$$
$$E[y_i] = g^{-1}(x_i^T \beta)$$
4. Link function g is a monotonic differentiable function (ensure inverse g^{-1} exists)
5. Estimate β by $\hat{\beta} = \text{argmax likelihood or argmax posterior probability}$
6. $E[y_i] = g^{-1}(x_i^T \beta)$

Linear regression: Natural link function g is the identity $\theta \in \mathbb{R}$

Poisson: " : " " is the log $\theta \in \mathbb{R}^{+*}$

Binomial: " : " " is logit function $\theta \in [0, 1]$ $g(\theta) = \log \left[\frac{\theta}{1-\theta} \right]$

Survival Analysis: " " will be log function

These proposed link functions relate to the function $b(\theta)$ defined for distributions in canonical form in exp family of distributions. Other link functions can be used

Survival Analysis	- Weibull, exponential	$y \in \mathbb{R}^+$
Linear Regression	- Normal	$y \in \mathbb{R}$
Poisson Regression	- Poisson	$y \in \mathbb{N}$
Logistic Regression	- Binomial	$y \in \{0, 1, \dots, n\}$

likelihood \equiv cost function corresponds to joint dist of all values given parameters needed

$P(y_1, y_2, \dots, y_n | \theta_1, \dots, \theta_n) = \prod_{i=1}^n P(y_i | \theta_i)$ independent of response variable y_i
only concerned with θ_i , other θ 's don't give info on y_i

$$\forall i \quad y_i \sim P_{y_i | \theta_i}(y_i | \theta_i)$$

$$\text{likelihood} = P(y_1, y_2, \dots, y_n | \theta_1, \dots, \theta_n) = \prod_{i=1}^n P_{y_i | \theta_i}(y_i | \theta_i)$$

$$\text{the log likelihood} = \sum_{i=1}^n \log(P_{y_i | \theta_i}(y_i | \theta_i))$$

$$P_{y_i | \theta_i}(y_i | \theta_i) = \exp[a(y_i) b(\theta_i) + c(\theta_i) + d(y_i)]$$

$$\log(L) = \sum_{i=1}^n [a(y_i) b(\theta_i) + c(\theta_i) + d(y_i)]$$

$$\text{each } \theta_i \rightarrow g^{-1}(x_i^T \beta) \quad \text{link function}$$

$$\hat{\beta} = \arg \max \log(L)$$

6

26/01/16 ALM2

LOGISTIC REGRESSION

Binomial dist: $P(y|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$ $y \in \{0, 1, \dots, n\}$

$$\begin{aligned}
 P(y|\theta) &= \exp \left[\log \left[\binom{n}{y} \theta^y (1-\theta)^{n-y} \right] \right] \\
 &= \exp \left[\log \left(\binom{n}{y} \right) + y \log(\theta) + (n-y) \log(1-\theta) \right] \\
 &= \exp \left[\underbrace{\log \left(\binom{n}{y} \right)}_{d(y)} + y \underbrace{\log(\theta)}_{a(y)} + (n-y) \underbrace{\log(1-\theta)}_{b(\theta)} \right]
 \end{aligned}$$

Value of θ that maximises $P(y|\theta)$?

$$\begin{aligned}
 \frac{d P(y|\theta)}{d(\theta)} &= \binom{n}{y} \left[y \theta^{y-1} (1-\theta)^{n-y} - \theta^y (n-y) (1-\theta)^{n-y-1} \right] \\
 &= \binom{n}{y} \left[\theta^{y-1} (1-\theta)^{n-y-1} \right] \left[y(1-\theta) - \theta(n-y) \right] \\
 y(1-\theta) - \theta(n-y) &= 0 \\
 y &= n\theta \quad \text{or} \quad \theta = y/n
 \end{aligned}$$

Expected Value

$$\begin{aligned}
 E[y] &= \sum_{y=0}^n y \binom{n}{y} \theta^y (1-\theta)^{n-y} \\
 &= \sum_{y=1}^n y \binom{n}{y} \theta^y (1-\theta)^{n-y} \quad \text{can't have } y=0 \text{ or else have 0.} \\
 y \binom{n}{y} &= \frac{y n!}{(n-y)! y!} = \frac{y \cdot n \cdot (n-1)!}{(n-1-(y-1))! y(y-1)!} = n \binom{n-1}{y-1} \\
 &= n\theta \sum_{y=1}^n \binom{n-1}{y-1} \theta^{y-1} (1-\theta)^{(n-1)-(y-1)} = n\theta \\
 &= \text{Bin}(\theta, n)
 \end{aligned}$$

$$y_1 \sim p(\theta, n_1)$$

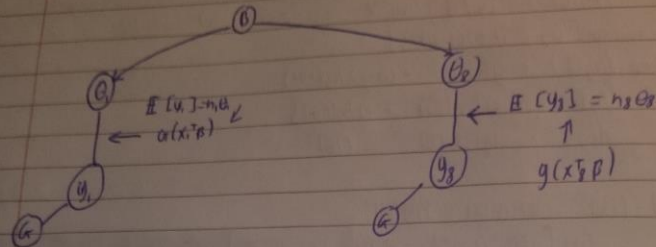
$$y_2 \sim p(\theta, n_2)$$

:

$$y_n \sim p(\theta, n_n)$$

26/01/16

Saturated model / Unconstrained GLM



$$\begin{aligned} L(\theta_1, \theta_2, \dots, \theta_k) &= P(y_1, \dots, y_k | \theta_1, \dots, \theta_k) \\ &= \prod_{i=1}^k P(y_i | \theta_i) \quad \text{independence} \\ &= \prod_{i=1}^k \binom{n_i}{y_i} \theta_i^{y_i} (1 - \theta_i)^{n_i - y_i} \end{aligned}$$

$$(\hat{\theta}_1, \dots, \hat{\theta}_k) = \text{argmax}_{\theta} L(\theta_1, \dots, \theta_k)$$

$$\hat{\theta}_1 = y_1/n_1, \dots, \hat{\theta}_k = y_k/n_k$$

Can plot $\hat{\theta}_i$ vs x_i

$$\theta \in [0, 1] \xrightarrow{g} x; p \in \mathbb{R} \text{ has to be invertible to get back to } g^{-1}$$

$$g(\theta) = \log \left[\frac{\theta}{1-\theta} \right] = z \quad z \in \mathbb{R}$$

 $[0, 1] \rightarrow (-\infty, +\infty)$ called logit function

$$g^{-1}(z) = \frac{e^z}{1 + e^z}$$

Can also use the probit function Φ

$$P(y | \beta, x) = \binom{n}{y} [g(\beta^T x)]^y [1 - g(\beta^T x)]^{n-y}$$

$$L(\beta) = \prod_{i=1}^N P(y_i | \beta, x_i)$$

$$= \prod_{i=1}^N \binom{n_i}{y_i} [g(\beta^T x_i)]^{y_i} [1 - g(\beta^T x_i)]^{n_i - y_i}$$

logit / probit etc

26/01/16 ALM2

3

Find β such that $dL(\beta)/d(\beta) = 0$

NOTE: you can add a $\beta_2 x^2$ term if you want to add another dimension to model

Dist	Links	Possible Models
Binomial	Probit	$\beta_0 + \beta_1 x$
	Logit	$\beta_0 + \beta_1 x + \beta_2 x^2$
		β_0
		$\beta_1 x$
		$\beta_1 x + \beta_2 x^2$
		$\beta_0 + \beta_2 x^2$ etc

- Use AIC/BIC to determine "best" model
- Can calculate deviance to see if model is good

GLM Example

g_i : Patient Survive? 1: No 0: yes

$f(x_i, y_i)$ $i=1, \dots, 40$

Response $y_i \sim \text{Binomial}$ $P(y_i|\theta) = \binom{n}{y} \theta^y (1-\theta)^{n-y}$

$$i=1 \quad P(y_1|\theta) = \binom{1}{y_1} \theta^{y_1} (1-\theta)^{1-y_1} \quad (\text{only 1 choice} \rightarrow \text{bernoulli})$$

$$= \theta^{y_1} (1-\theta)^{1-y_1}$$

$$i=2 \quad = \theta^{y_2} (1-\theta)^{1-y_2}$$

$$i=40 \quad = \theta^{y_{40}} (1-\theta)^{1-y_{40}}$$

$$E[y] = n\theta = \theta \quad n=1 (\text{bernoulli})$$

$$\frac{dL(y|\theta)}{d(\theta)} = 0 \quad \Rightarrow \theta = y/n = y$$

$$L(\theta_1, \dots, \theta_{40}) = \prod_{i=1}^{40} \theta_i^{y_i} (1-\theta_i)^{1-y_i} \quad \dim(\vec{\theta}) = 40$$

$$(\hat{\theta}_1, \dots, \hat{\theta}_{40}) = \operatorname{argmax} L(\theta_1, \dots, \theta_{40})$$

$$\text{link function} = b(\theta) = g(\theta) = \log\left[\frac{\theta}{1-\theta}\right] \rightarrow \logit$$

$$\text{GLM: } L(\beta) = \prod_{i=1}^{40} g^{-1}(\beta_0 + \beta_1 x_i)^{y_i} [1 - g^{-1}(\beta_0 + \beta_1 x_i)]^{1-y_i}$$

$$\dim(\beta) = 2$$