

# Logit Regression.

LA-1

Categorical variables.

→ Why Linear Reg. Not possible  
Response Not on a ratio scale.

Error terms not normally distributed.

Linear Reg. can generate any real value.

but categorical var can only take limited nos of discrete values.

→ Fundamental difference of Logit with L.R.  
L.R. → Exp. value of the dep. var = Linear combination of indep vars. and their corresp. parameters.

Generalized linear models. equate

the linear combination of indep vars. to some function of the probability of a given outcome on the dep. var.

Logit. Regr.: } That function is Logit transform. Natural log. of the odds that some event will occur.

Parameter estimation in Linear Reg.

- Minimizing sum of squares of dev. of predicted values from observed values.
- Solution involves solving system of  $n$  equations.

The Model → Binomial Logit Regression.

Consider  $i$ th Row → has distinct Combo.

of values

$N \rightarrow$  total Nos of populations

of indep. variables.

$\begin{bmatrix} Y_i \end{bmatrix} \xrightarrow{N} \text{Nos of successes of } z \text{ for population } i$

$\begin{bmatrix} Y_i \end{bmatrix} \rightarrow \text{Observed counts of the Nos of Success For each population.}$

$\begin{bmatrix} \pi_i \end{bmatrix} \rightarrow \text{Prob of success. For any given obs in } i\text{th population}$

The Logit Model ] Equates the logit transform (Log odds) of the Prob. of success.

$$\text{Log} \left( \frac{\pi_i}{1 - \pi_i} \right) = \sum_{k=0}^K x_{ik} \beta_k, \quad i = 1, 2, \dots, N$$

Getting an expression For  $\pi_i = P(x)$

$$\frac{P(x)}{1 - P(x)} = \exp \left( \sum_{j=0}^K x_j \beta_j \right) = \prod_{j=0}^K \exp(x_j \beta_j)$$

$$\text{Then } \frac{P(x)}{1 - P(x)} = \exp(Z)$$

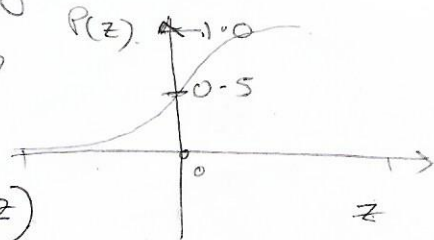
$$\text{Let } Z = \prod_{j=0}^K x_j \beta_j$$

Maps a Real Line to (0,1)

$$\text{or } P(x) = \frac{\exp(Z)}{1 + \exp(Z)} = P(Z)$$

Sigmoid Function.

$$\text{Thus } P(Z) = \frac{\exp(Z)}{1 + \exp(Z)} = \frac{1}{1 + \exp(-Z)}$$



$$p'(z) = p(z)(1-p(z)) \rightarrow \text{Root} \quad p(z) = \frac{1}{1-e^{-z}} = (1-e^{-z})^{-1}$$

$$\therefore p'(z) = -1(1-e^{-z})^{-2}(e^{-z})$$

$$\text{now } p(z) = \frac{1}{1-e^{-z}} \quad = \frac{-e^{-z}}{(1-e^{-z})(1-e^{-z})}$$

$$\text{and } 1-p(z) = 1 - \frac{1}{1-e^{-z}}$$

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

$$\text{Thus. } [p'(z) = p(z)(1-p(z))]$$

Note  $p'(\beta) = p(z)(1-p(z)) \cdot (z')$  where  $z'$  is gradient taken wrt  $\beta$ .

Maximum Likelihood Estimation.

$L(X|P) = \prod$  Product of predicted probabilities of the  $N$  individual observations.

$$L(X|P) = \prod_{i=1}^n p(x_i) \prod_{i=1}^n (1-p(x_i))$$

$y_i=1 \quad y_i=0$

$X_{(K+1, N)}$  Each column corresponds to an observation.  
First row is 1.

$y$   $N$  dim vector of Responses.

$(X, y) \rightarrow$  is a set of observations.



Log likelihood.

[LR=4]  
[ - . - ]

$$\mathcal{L}(x|p) = \sum_{i=1, y_i=1}^N \log p(x_i) + \sum_{i=0, y_i=0}^N \log (1-p(x_i))$$

we want to maximize Log likelihood.

$$\nabla_b \mathcal{L} = \sum_{\substack{i=0 \\ y_i=1}}^N \frac{p_i'}{p_i} x_i - \sum_{\substack{i=0 \\ y_i=0}}^N \frac{p_i' (x_i)}{(1-p_i)} x_i$$

Now  $p_i' \equiv p(1-p)$ .

Thus we actually have.

$$\nabla_b \mathcal{L} = \sum_{\substack{i=1 \\ y_i=1}}^N \frac{p_i (1-p_i)}{p_i} x_i - \sum_{\substack{i=1 \\ y_i=0}}^N \frac{p_i (1-p_i)}{1-p_i} x_i$$

$$\text{Thus } \nabla_b \mathcal{L} = \sum_{\substack{i=1 \\ y_i=1}}^N (1-p_i) x_i - \sum_{\substack{i=1 \\ y_i=0}}^N p_i x_i$$

$$= \sum_{i=1}^N [y_i (1-p_i) - (1-y_i) p_i] x_i$$

after cancelling terms.

$$= \sum_{i=1}^N [y_i x_i - p_i x_i] = 0$$

Notes:

Equations to be solved are in terms of probabilities ( $p$ ). Not in terms of  $b$ .

$\therefore$  Logit Models are coordinate free.

f(Solving For coeff.)

[4R-5]

Say for a vector valued func.  $f(b)_{opt} = 0$

Taylor expansion around initial guess.

$$f(b_0 + \Delta) \approx f(b_0) + f'(b_0) \Delta \quad f' \rightarrow \text{Jacobian of first derivatives of } f \text{ wrt } b.$$

Now setting LHS to zero, we have.

$$\Delta_0 = [f'(b_0)]^{-1} f(b_0)$$

Now update estimate for  $b$ :  $b_1 = b_0 + \Delta_0$  iterate till convergence.

Now  $H = \frac{\partial}{\partial b} \nabla_b L$

Multi Variable Taylor's Series.

$$f(x, y) = f(a, b) + f_x(a, b)(x-a) + f_y(a, b)(y-b) + \frac{1}{2!} [f_{xx}(a, b)(x-a)^2 + 2f_{xy}(a, b)(x-a)(y-b) + f_{yy}(a, b)(y-b)^2] + \dots$$

Now we already have:

$$\nabla_b L = \sum_{i=1}^N [y_i x_i - p_i x_i]$$

Now

$$H = \frac{\partial}{\partial b} \nabla_b L = - \sum_{i=1}^N x_i \nabla_b p_i = - \sum_{i=1}^N x_i p_i (1-p_i) x_i^T$$

$$\text{or } H = X W X^T$$

Thus we have the following.

[LR-5]

$$\Delta_k = (X W_k X^T)^{-1} X (y - p_k)$$

where we have the following.

- >  $W$  is a diagonal matrix of the derivatives  $p_i'$
- > The  $i$ th. col of  $X$  corresponds to  $i$ th observation.

Compare with LINEAR Regression.

$$y = X^T b$$

$$Xy = XX^T b$$

$$b = (XX^T)^{-1} Xy$$

Comparing the two we find. that at each iteration  $\Delta$  is the solution of a weighted Least squares.

Notes:

Coeff. tends to infinity  $\rightarrow$  Sign that a input is perfectly correlated.

Large. Coeff. Magnitudes:  $\rightarrow$  indication of correlated inputs.

[The AIC Value to Compare Models.] Akaike Information Criterion (AIC)

$$AIC = 2K - 2 \ln(\hat{L})$$

Maximized value of the likelihood function

Nos of Estimated Parameters

Lower the AIC Value better the model

Notes:]

- 1] It Supports better Fit
- 2] Penalizes using more Parameters.



# Classifying using Linear Regression.

- Train Model to Predict  $[0, 1]$
- $> 0.5 \rightarrow 1$   $< 0.5 = 0$
- Cons] → Raw output from LR → can Predict value outside this Range
- LR → designed to Min Sq Err. against line of best fit.
- ↓
- Thus decision boundary highly sensitive to influential observations.
- Regression Line shifts

## Logit Function

Sigmoid

$$f(x) = \frac{1}{1 + e^{-x}} \quad \text{output } [0, 1] \text{ always.}$$

## Deviance Computation

↳ Used in lieu of Sum of Squares Computation.

$$D = -2 \ln \frac{\text{likelihood of fitted model}}{\text{likelihood of Saturated model}}$$

D approx follows Chi square distribution.

When Saturated model not avail.

$$D = -2 \ln (\text{likelihood of fitted model})$$

Null Deviance → Model with just intercept vs Saturated.

$$D_{\text{null}} = -2 \ln \left( \frac{\text{likelihood of Null model}}{\text{likelihood of Saturated model}} \right)$$

## Model Deviance

$$D_{\text{fitted}} = -2 \ln \left( \frac{\text{likelihood of fitted model}}{\text{likelihood of Saturated model}} \right)$$

## Good Fit

For Model to have good fit.

Model Deviance Should be Significantly Smaller than Null Deviance.

## Likelihood Function . For Logeshtics Regression

Notes  $\rightarrow$  For Each training data point we have a vector of features  $x_i$ , and an observed class  $y_i$

Prob of that class was  $p$  if  $y_i = 1$   
or  $1-p$  if  $y_i = 0$

Thus one can write

the likelihood as.

$$L(\beta_0, \beta) = \prod p(x_i)^{y_i} (1-p(x_i))^{1-y_i}$$

The log likelihood is then.

$$l(\beta_0, \beta) = \sum_{i=1}^n y_i \log p(x_i) + (1-y_i) \log (1-p(x_i))$$

## [Goodness of Fit test]

Pseudo  $R^2 =$

$$\frac{D_{\text{NULL}} - D_{\text{fitted}}}{D_{\text{NULL}}}$$

$$= 1 - \frac{D_{\text{fit}}}{D_{\text{NULL}}}$$



Logit Regression  
(Log Likelihood Max)

Logit R.  
(9)

$$P(z) = \frac{e^z}{1+e^z} = (e^z)(1+e^z)^{-1} \quad \left[ \begin{array}{l} \text{Note} \\ 1 - P(z) = 1 - \frac{e^z}{1+e^z} = \frac{1}{1+e^z} \end{array} \right]$$

$$\begin{aligned} P'(z) &= (e^z)(1+e^z)^{-1} + (-1)(e^z)(1+e^z)^{-2} e^z \\ &= \frac{e^z}{1+e^z} - \frac{(e^z)^2}{(1+e^z)^2} = \frac{e^z(1+e^z)}{(1+e^z)^2} - \frac{(e^z)^2}{(1+e^z)^2} \\ &= \frac{e^z}{(1+e^z)} \cdot \frac{1}{(1+e^z)} \end{aligned}$$

Thus  $P'(z) = P(z)(1-P(z))$

Log likelihood  $\rightarrow \mathcal{L}(x/p) = \sum_{i=1, y_i=1}^N \text{Log } P(x_i) + \sum_{i=0, y_i=0}^N \text{Log } (1-P(x_i))$

$$\nabla_b \mathcal{L} = \sum_{\substack{i=0 \\ y_i=1}}^N \frac{p_i'}{p_i} x_i - \sum_{\substack{i=0 \\ y_i=0}}^N \frac{p_i'}{1-p_i} x_i$$

$$\begin{aligned} &\Rightarrow \sum_{\substack{i=1 \\ y_i=1}} (1-p_i) x_i - \sum_{\substack{i=1 \\ y_i=0}} p_i x_i \\ &= \sum_{i=1}^N [y_i(1-p_i) - (1-y_i)p_i] x_i \end{aligned}$$

Thus we arrive at the set of Simultaneous equations that are true at the optimum.

$$\sum_{i=1}^N y_i x_i - p_i x_i = 0$$

How to solve for coefficients?

Solving for coeff:  $\rightarrow$

[Logit R]

Suppose you have a vector valued function.  
 $f: y = f(b)$  you want the value  $b_{opt}$ .

(10)

Such that  $f(b)_{opt} = 0$ , Assume we start with initial guess  $(b_0)$

$$\underbrace{f(b_0 + \Delta)}_{0''} \approx f(b_0) + f'(b_0) \Delta$$

$$\Delta_0 = - \frac{f'(b_0)}{f''(b_0)}$$

$$b_1 = b_0 + \Delta_0$$

In our case.  $f = \nabla_b L = 0$  ] Now  $\nabla_b L = \sum y_i x_i - p_i x_i = 0$   
 $H = \frac{\partial}{\partial b} \nabla_b L$  ]  $H = \frac{\partial}{\partial b} (\nabla_b L) = - \sum_{i=1}^N x_i (\nabla_b p_i)$   
 $= - \sum_{i=1}^N x_i p_i (1-p_i) x_i^T$

$$\text{thus } \Delta_k = (X W_k X^T)^{-1} X (y - p_k)$$

Compare with Linear Regression.

$$y = X^T b$$

$$X y = X X^T b$$

$$b = (X X^T)^{-1} X y$$

$$= X W X^T$$

L.R.  $\rightarrow$

Like a weighted  
Least squares problem