

Zaki (2026)

624  
635

①

## logistic Regression

24.1 Binary

623

24.1.1 MLE

626

24.2 Multiclass

630

24.2.1 MCE

632

Aim of log reg is to predict prob of response variable based on the indep variables

log reg = classification Tech that  
g for a given point  $x_i \in \mathbb{R}^d$  preds  
 $P(C_j | x_i)$  for each class in  
the domain of  $y$

log reg: Indep vars  $x_1, x_2, \dots, x_n$   
Bin or Bernoulli response  $y \in \{0, 1\}$

~~At~~ like linear Reg, Aug data to add

$x_0$  value = 1 for bias weight

$$\tilde{x}_i = (1, x_1, x_2, \dots, x_d)^T \in \mathbb{R}^{d+1}$$

(2)

In the Binary log reg there are only two outputs: so the prob mass function for  $\tilde{X} = \tilde{x}$  is:

$$P(Y=1 | \tilde{X} = \tilde{x}) = \pi(\tilde{x})$$

means prob function

$$P(Y=0 | \tilde{X} = \tilde{x}) = 1 - \pi(\tilde{x})$$

$$\text{So, } \pi(\tilde{x}) = Y=1 \text{ given } \tilde{X} = \tilde{x}$$

Also denotes Expected value of  $Y$  given  $\tilde{X} = \tilde{x}$

thus, in logistic regression the goal is to predict the prob

this is in conflict to say Bayes where the prob is answered & classif matter most

Note:

$$P(Y=1 | \tilde{X} = \tilde{x}) \text{ also = Exp Val of } Y \text{ give } \tilde{X} = \tilde{x}$$

(3)

Since we want a probability it is not correct to use the linear reg model

This because we need output to be between 0-1

the sigmoid function achieves this

$$\theta(z) = \frac{1}{1 + \exp\{-z\}} = \frac{\exp\{z\}}{1 + \exp\{z\}}$$

logistic reg model is defined as follows:

$$P(Y=1 | \tilde{x} = \tilde{x}) = \pi(\tilde{x}) = \theta(f(\tilde{x})) = \theta(\tilde{w}^T \tilde{x}) = \frac{\exp\{\tilde{w}^T \tilde{x}\}}{1 + \exp\{\tilde{w}^T \tilde{x}\}}$$

$$P(Y=0 | \tilde{x} = \tilde{x}) = 1 - P(Y=1 | \tilde{x} = \tilde{x}) = \theta(-\tilde{w}^T \tilde{x}) = \frac{1}{1 + \exp\{\tilde{w}^T \tilde{x}\}}$$

(4)

here the probs for  $y=1$  &  $y=0$  are demonstrated using different versions of the sigmoid functions

the  $y=1$  uses the form  $\frac{\exp(w^T x)}{1 + \exp\{w^T x\}}$

this is because it shows how the prob of  $y=1$  is related to the expo of the linear comb of  $w^T x$

Experiments that as the ~~expo~~ C increase so does the prob

Why?

- The term is in the numerator so over -
- all value of the fraction must increase as it does
- ~~term~~ remember it is expo so as  $w^T x$  increase output does by allot

(5)

$y = \Theta$  uses the form  $\frac{1}{1 + \exp(w^T x)}$

this emphasizes that as  $w^T x$  increases, the prob is related to the inverse of the expo

As  $w^T x$  increases the prob decreases

Denominator grows as  $w^T x$  does

NOTE: both forms are the same and give the same result, this is just used for easy reading & interpretation

Combining the two outcome cases  
 $(y=0, y=1)$  the full logistic reg  
 Model is given as:

$$P(Y = \tilde{X} = \tilde{x}) = \Theta(\tilde{\omega}^T \tilde{x}) \cdot \Theta(-\tilde{\omega}^T \tilde{x})^{1-y}$$

Note  $1 - \Theta(z)$  same as  $\Theta(-z)$

$$\text{So } P(Y | \tilde{X} = \tilde{x}) = \Theta(\tilde{\omega}^T \tilde{x}) = Y=1$$

$$P(Y | \tilde{X} = \tilde{x}) = \Theta(-\tilde{\omega}^T \tilde{x}) = Y=0$$

# Log Odd Ratio

for occurrence  $Y=1$ :

$\leftarrow$  just odds, not log

$$\text{odds}(Y=1 | \tilde{X} = \tilde{x}) = \frac{P(Y=1 | \tilde{X} = \tilde{x})}{P(Y=0 | \tilde{X} = \tilde{x})}$$

= Ratio of Probabilities

= General rule = the "Happen" class on top

$$= \frac{\exp\{\tilde{\omega}^T \tilde{x}\}}{1 + \exp\{\tilde{\omega}^T \tilde{x}\}} \cdot \left(1 + \exp\{\tilde{\omega}^T \tilde{x}\}\right)$$

$$= \exp\{\tilde{\omega}^T \tilde{x}\}$$

Next, we want to transform this into the log odds ratio. Why?

- Odds are uncapped: 0 to  $\infty$ . Large prob are hard to deal with

- Log odds transforms to Neg inf to Pos inf

- Makes it easier to compare coeffs

- A positive log odds = increased odds of  $Y=1$

- A negative log odds = decreases odds of  $Y=1$

(8)

- Creates a linear odds ratio w/ preds
- ratio of probs & sigmoid allow for wif & exponential values
- Allows for symmetry of numbers used.  
w/ odds .2 &  $\frac{1}{2}$  represent same  
magnitude of change in diff directions  
(half & double) but are very diff  
numerically. log odds creates a .0  
bounds & symmetry:  $2 = 0.7$ ,  $\frac{1}{2} = -0.7$

the logarithm of the odds ratio:

$$\ln(\text{odds}(Y=1 | \tilde{x} = \tilde{x})) = \ln \left( \frac{P(Y=1 | \tilde{x} = \tilde{x})}{P(Y=0 | \tilde{x} = \tilde{x})} \right)$$

$$= \ln(\exp\{\tilde{w}^T \tilde{x}\}) = \underline{\tilde{w}^T \tilde{x}}$$

This is just the output of the linear Model

- log odds are linearly related to predictors
- coeffs ( $w$ ) rep change in log odds for one-unit A in pred
- log reg is modelling log-odds as linear preds of the predictors

(9)

the log odd ratio function is called logit

$$\text{logit}(z) = \ln\left(\frac{z}{1-z}\right)$$

where  $z =$   
odd ratio  
calced by  
 $\frac{P}{1-P}$

This is the inverse of the logistic function

$$\underbrace{\ln(\text{odds}(Y=1 | \tilde{X}=\tilde{x}))}_{\text{Same as } \frac{1}{1+\text{exp}(w^T x)}} = \text{logit}(P(Y=1 | \tilde{X}=\tilde{x}))$$

$$\text{logistic: } \frac{1}{1+\text{exp}(w^T x)}$$

- takes in real number as input
- outputs value between 0 & 1
- Maps real number to Prob

Logit:

- takes Prob or odd ratio as input
- Produces log odds
- Maps Prob to real number

Hence the two reverse each other

the logistic regression model is therefore based on the assumption that the log odds ratio for  $Y=1$  given  $\tilde{X}=\tilde{x}$  is a linear function of the independent variables.

## 24.1.1 Maximum likelihood Estimation

let  $D$  be train w/  $n$  points  $x_i$ , & labels

let  $\tilde{w} = (w_0, w_1, w_2, \dots, w_n)^T$  be the  
augment weight vector for est  $\tilde{w}$

We will uses Max likelihood  
approach to learn the weight  
vector

likelihood = prob of observing data  
given the est parameters  $\tilde{w}$

~~Then~~ likelihood of observed resp  
is given as:  $L(\tilde{w}) = P(Y = \tilde{w}) =$

$$= \prod_{i=1}^n P(y_i | \tilde{x}_i) = \prod_{i=1}^n \theta(\tilde{w}^T \tilde{x}_i)^{y_i} \cdot \theta(-\tilde{w}^T \tilde{x}_i)^{1-y_i}$$

However, instead of trying to max the likelihood, we can max the log of likelihood to convert the product into a summation

$$\ln(L(\tilde{w})) = \sum y_i \cdot \ln(\Theta(\tilde{w}^T \tilde{x}_i)) + (1-y_i) \cdot \ln(\Theta(-\tilde{w}^T \tilde{x}))$$

the neg of the log-likelihood can be considered an error func:

Cross-entropy error function

$$E(\tilde{w}) = -\ln(L(\tilde{w})) = -\sum y_i \cdot \ln(\Theta(\tilde{w}^T \tilde{x}_i)) - (1-y_i) \cdot \ln(1-\Theta(\tilde{w}^T \tilde{x}_i))$$

the task of maxing the log-like is there the same as min the cross entropy error

Typically, to solve the optimal weight vector  $\hat{w}$  we would diff the log-like funct w/ resp to  $\tilde{w}$  & set to 0 and solve for  $\tilde{w}$

However, there results in no-closed form to solve for  $\tilde{w}$ . That is not  $\tilde{w}$  term remains

Instead use + iterative gradient Ascent to comp optimal value as log-likli function is concave

This method relies on the gradient of the log-lik which is obtained by partial Derivative w/ respect to  $\tilde{w}$

$$\text{Change in } \hat{w} = \nabla \hat{w} = \frac{\partial}{\partial \tilde{w}} \left\{ \ln(L(\tilde{w})) \right\} =$$

$$\frac{\partial}{\partial \tilde{w}} \left\{ \sum_i y_i \cdot \ln(\theta(z_i)) + (1 - y_i) \cdot \ln(\theta(-z_i)) \right\}$$

$z_i = \hat{w}^T \tilde{x}$  this is just the log-likli

To obtain  $\nabla(\tilde{w})$  we need to use the chain rule to obtain deriv of  $\ln(\theta(z_i))$  w.r.t respect to  $\tilde{w}$

$$\text{for } y=1 \rightarrow \theta(-z_i) \cdot x_i$$

$$y=2 \rightarrow -\theta(z_i) \cdot \tilde{x}_i$$

replace  $\ln(\theta(z_i))$  for respective

$$\nabla(\tilde{w}) = \sum_{i=1}^n w_i \cdot \theta(-z_i) \cdot \tilde{x}_i - (1-w_i) \cdot \theta(z_i) \cdot \tilde{x}_i$$

$$= \sum_{i=1}^n (\theta(-z_i) + \theta(z_i)) \cdot x_i - \theta(z_i) \cdot x_i$$

$$= \sum_{i=1}^n (y_i - \theta(z_i)) \cdot \tilde{x}_i \quad \text{since } y_i = \theta(z_i)$$

$$= \sum_{i=1}^n (y_i - \theta(\tilde{w}^\top \tilde{x})) \cdot \tilde{x}_i \quad \nabla(\tilde{w})$$

the sum of the class minus the linear output times the data point

from here the gradient ascent algorithm starts for some initial estimate of  $\tilde{w} = \hat{w}^0$ ,

At each step  $t$  the method moves in the direction of the steepest ascent which is given by the gradient vector  $\nabla \tilde{w}$

Note:  $\tilde{w}^t$  = est/start weight vector

$\Delta \tilde{w}^t$  =  $w^t$  plugged into part b(c)

$\Delta \tilde{w}^t$  is the direction & size to update

if  $\tilde{w}^t$  is current we obtain the next by:

$$\tilde{w}^{t+1} = \tilde{w}^t + \eta \cdot \nabla(\tilde{w}^t)$$

$\eta > 0$  is user specific param called learning rate

if too large estimates will vary too much

too small & converges will be (very) slow

@ optimal  $\tilde{w}$ ,  $\nabla(\tilde{w}) = 0$

(16)

It is important to note / remember that each partial derivative measures how sensitive the function is to changes in that specific variable ( $w$ )

A larger Partial Deriv = a more rapid movement in a direction

Note, when updating we are updating the whole weight vector

the gradient vector  $\nabla w$  is the output of the partial derivatives for each weight

~~The gradient vector is updated after each partial derivatives stage~~

~~Method of Conjugate Gradient~~

# Stochastic Gradient Ascent

the previous gradient Ascent calc  
Computes by considering all data  
Points  $X$ ,

It is also called Batch GA

for large database, faster to compute  
gradient by considering one randomly  
choose point

the weight vector is updated  
at each partial Derivatives  
stage

= SGA for compute optimal  $\tilde{w}$

$$\nabla(\tilde{w}_i, \tilde{x}_i) = (y_i - \theta(\tilde{w}^T \tilde{x}_i)) \cdot \tilde{x}_i$$

SGA = much Quicker & cheaper comp  
to Batch

## Predictions

Once the model has been trained  
i.e. we know the optimal  $\tilde{w}$   
weight vector

we predict using

$$\hat{y} = \begin{cases} 1 & \text{if } \theta(\tilde{w}^T \tilde{x}) \geq 0.5 \\ 0 & \text{if } \theta(\tilde{w}^T \tilde{x}) < 0.5 \end{cases}$$