

Logistic Regression: a few facts.

Facts

- A regression model for categorical endpoints / responses (here: binary) used to analyze the impact (magnitude, direction, inference) of a predictor variables on the probability of occurring an event.
- As it's a regression, it predicts a numerical outcome with $E(Y|x=x)$.
For the Bernoulli's distribution (or binomial with $k=1$) - it's a probability of success (or log-odds, depending on formulation).
more precisely, it's a direct probability estimate.
- Logistic regression itself does NOT produce labels. It's a direct probability estimator.
To give classes, it needs additional post-fit activity:
applying a decision rule to the predicted probability.

This makes the logistic classifier.

(Further explained in detail)

- o LR is fit by the maximum likelihood
(it's NOT just OLS with logit transformation)
- o LR is a part of the Generalized linear model
(like linear, Poisson, gamma, beta, etc.).

Logistic regression, invented and further developed in terms of the Generalized Linear Model and Generalized Additive Model by Cox, Nelder, Hastie, Tibshirani and others, served exactly to address regression problems with categorical outcomes (via maximum likelihood estimation).

Generalized Linear models.

In regression example, we had $y|\alpha; \theta \sim N(\mu, \sigma^2)$ and in the classification one, $y|\alpha; \theta \sim \text{Bernoulli}(\phi)$ for some appropriate definitions of μ and ϕ as function of α and θ .

Both of these cases of methods are special cases of a broader family of models, called Generalized Linear models (GLMs)

Exponential Family.

To work our way up to GLMs, we will begin by defining exponential family distribution.

A class of distribution is said to be in the exponential family if it can be written in the form

$$\{ p(y; \eta) = b(y) \exp(\eta^T T(y) - a(\eta)) \}$$

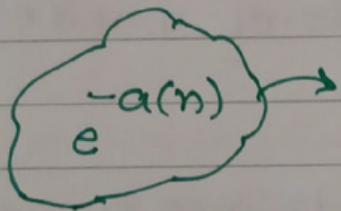
Here,

η = natural parameter, also called the canonical parameter of the distribution

~~abuse word~~ beginning

$T(y) = \text{sufficient statistic}$

$a(n) = \text{log-partition function}$



plays the role of a normalization constant,

that makes sure the distribution $p(y; n)$ sums/integrates over y to 1.

A fixed choice of T , a and b defines a family of distributions that is parameterized by n , we get then different distributions within this family.

Bernoulli Distribution

The Bernoulli distribution with mean(ϕ);
 $\text{Bernoulli}(\phi)$,

specifies a distribution over $y \in \{0, 1\}$, so that

$$\{ p(y=1; \phi) = \phi ; p(y=0; \phi) = 1-\phi \}$$

as we vary ϕ , we obtain Bernoulli distributions with different means.

Now, we will show that the class of Bernoulli distributions obtained by varying ϕ , is in Exponential family.

We write the Bernoulli distribution as:

$$\begin{aligned} p(y; \phi) &= \phi^y (1-\phi)^{1-y} \\ &= \exp(y \cdot \log \phi + (1-y) \cdot \log(1-\phi)) \\ &= \exp\left(\left(\log\left(\frac{\phi}{1-\phi}\right)\right)y + \log(1-\phi)\right) \end{aligned}$$

Thus,

$$\text{natural parameter } (\eta) = \log\left(\frac{\phi}{1-\phi}\right)$$

$$\begin{aligned} T(y) &= y \\ a(\eta) &= -\log(1-\phi) \\ &= \log(1+e^{-\eta}) \end{aligned}$$

$$b(y) = 1$$

This shows that the Bernoulli distribution can be written in the form of Exponential family, using an appropriate choice of T , a and b .

Similarly, we can also show for Gaussian distribution.

There're many other distributions that are members of the exponential family.

multinomial distribution, poisson, gamma, exponential distribution and many more are the members of exponential family.

Exponential Family \neq Exponential distribution

Constructing GLMs

To derive a GLMs, we will make the following three assumptions about the conditional distribution of y given x and about the model:

1. $y|x; \theta \sim \text{Exponential Family}(\eta)$

i.e., given x and θ , the distribution of y follows some exponential family distribution, with parameter η .

2. Given α , our goal is to predict the expected value of $T(y)$ given α .

Example: $T(y) = y$, so this means we would like the prediction $h(\alpha)$ output by our hypothesis h to satisfy $h(\alpha) \approx E[y|\alpha]$.

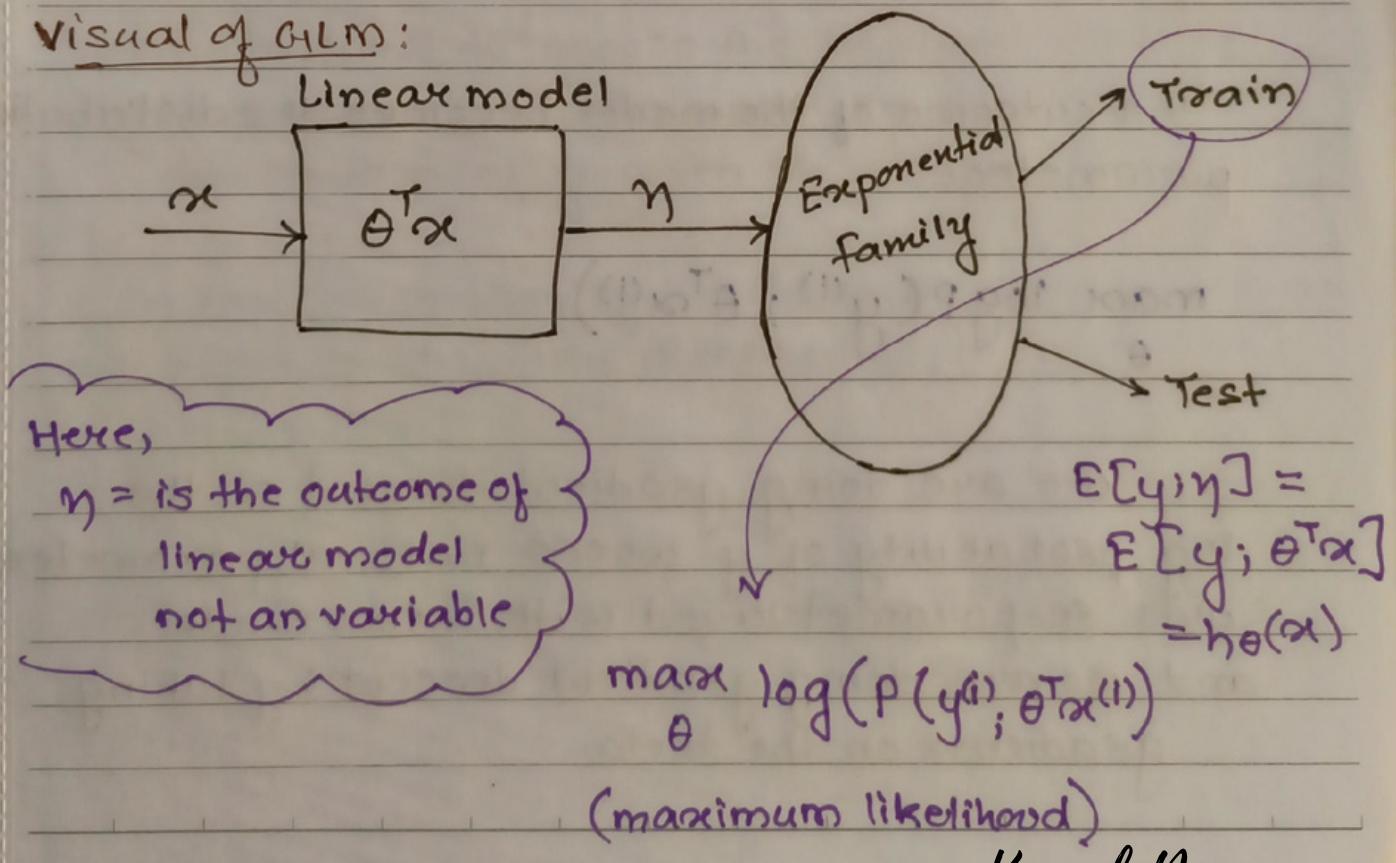
output: $E[y|\alpha; \theta]$

$$h_\theta(\alpha) = E[y|\alpha; \theta]$$

3. The natural parameter η and the inputs α are related linearly: $\eta = \theta^T \alpha$

(or, if η is vector-valued, then $\eta_i = \theta_i^T \alpha$)

Visual of GLM:



we are training θ to predict the parameter of the exponential family distribution whose mean is the prediction that we're gonna make for y .

How do we train this model?

So, in this model, the parameter that we are learning by doing gradient descent are $\theta^T \alpha$ parameters.

so we are not learning any parameters in the exponential family, we are learning θ that's part of the model and not part of the distribution.

The outcome of the model becomes the distribution parameters.

$$\max_{\theta} \log P(y^{(i)}; \theta^T \alpha^{(i)})$$

we are doing gradient descent on the log probability of y where natural parameter was re-parameterized with linear model and we are doing gradient descent by taking gradients on the theta.

GLM training:

learning update rule:

$$\theta_j := \theta_j + \alpha (y^{(i)} - h_{\theta}(x^{(i)})) x_j^{(i)}$$

This remains same for all distributions.

Logistic Regression

Here, we are interested in binary classification
so, $y \in \{0, 1\}$

Given that y is binary-valued, it therefore
seems natural to choose the Bernoulli family
of distributions to model the conditional
distribution of y given x .

In our formulation of the Bernoulli distribution as
an exponential family distribution,
we had

$$\phi = \frac{1}{1 + e^{-\eta}}$$

if $y|x; \theta \sim \text{Bernoulli}(\phi)$, then $E[y|x; \theta] = \phi$

$$h_{\theta}(\alpha) = E[y|\alpha; \theta]$$

$$\begin{aligned} &= \phi \\ &= \frac{1}{1 + e^{-\eta}} \\ &= \frac{1}{1 + e^{-\theta^T \alpha}} \end{aligned}$$

So, this gives us hypothesis functions of the form $h_{\theta}(\alpha) = \frac{1}{1 + e^{-\theta^T \alpha}}$

∴ once we assume that y conditioned on α is Bernoulli, it arises as a consequence of the defn of GLMs and exponential family distributions.

more terminology,

the function g giving the distribution's mean as a function of the natural parameter

$$g(\eta) = E[\tau(y); \eta]$$

is called the canonical response function.

Its inverse, g^{-1} , is called the canonical link function.

Thus, the canonical response function for the Bernoulli is the logistic function.

when the $p(y|\alpha)$ are assumed to be....

Poisson
Negative binomial

Poisson regression
Negative binomial regression

Bernoulli
Normal
Multinomial

Logistic regression
Classical regression
Multinomial

Beta
Lognormal
...

Beta regression
Lognormal regression

The intent is simply to show that different assumed families of conditional distributions $p(y|\alpha)$ give rise to different types of regression models.

To conclude,

Logistic regression is a special case of the generalized linear regression model (GLM).

In ML, it is used mainly for classification by taking the predicted probability and applying a decision rule to it.

This makes the logistic classifier.

Illustration →

Step 01

Logistic Regression

takes the training data and produces the training data with the model coefficients

that will be used for the prediction.

used the GLM:

Step 2

use the LR model
to predict p

$$p(\text{LR model}, X) = 0.76$$

Step 3

apply a decision rule to classify

p

If $p > 50\%$, label "B"

new
input
data

Logistic classifier based
on the logistic regression model.

Product:
classified label:
"B"

The logistic regression model itself simply models probability of output in terms of input and does not perform statistical classification.
(it is not a classifier),

though it can be used to make a classifier, for instance by choosing a cutoff value and classifying inputs with probability greater than the cutoff as one class, below the cutoff as the other ;

this is a common way to make a binary classifier.

Logistic Regression \neq
logistic classifier

How is it used for segregation purposes?

It assesses how many a certain set of independent variables affects the odd-ratio or % of success (depending on formulation) through the assessment of the

- o marginal and interaction effects:

dissection, magnitude, inference (p-value, confidence interval)

- o simple effects (for categorical covariates):

magnitude of different ratio + inference (p, CI)

.. and, when the model is followed by ANOVA-like type analysis, here called the "analysis of deviance", it gives us also the main and their intersection effects,

eg. "does gender, does gender and treatment affect the % of successes?"

It's done either via likelihood Ratio tests or wald's joint analysis of the model coefficients.

The analysis employs the EM-means (estimated marginal means), which again - is a typical segregation outcome.