# MLA
## LOGISTIC REGRESSION

- Parametric classification technique making distributional assumptions are made to allow for probabilistic quantification of class assignment

## EXAMPLE
- Resting pulse depends on patient smoking and their weight
- Data 92 old, 1=low, 0=high, Smoke 1=yes 0=no and weight in lb

- $P(Low)$ probability of low pulse rate
- To ensure $P(low)$ fall within $[0,1]$ we consider following model:

$$Logit(P(Low)) = Log\left(\frac{P(Low)}{1-P(Low)}\right) = \alpha + \beta_s Smoke + \beta_w Weight$$

or $\dfrac{P(low)}{1-P(low)} = exp(\alpha + \beta_s Smoke + \beta_w Weight)$

or $P(low) = \dfrac{exp(\alpha + \beta_s Smoke + \beta_w Weight)}{1 + exp(\alpha + \beta_s Smoke + \beta_w Weight)}$

## VARIANTS
- Logistic model may not be only choice that ensures $P(low) \in [0,1]$
- Can use probit model:

$$\Phi^{-1}(P(Low)) = \alpha + \beta_s Smoke + \beta_w Weight$$

Here $\Phi^{-1}(z) = w$ if probability of a normal $(0,1)$ random variable $v$ is less than or equal to $z$, it is the inverse CDF of $N(0,1)$

- Gompit / Complementary log-log: $-Log[-log(P(Low))] = \alpha + \beta_s Smoke + \beta_w Weight$
- Known as LINK FUNCTIONS

- More evidence required in logistic model than the probit model to increase probability of assignment to that particular class

## PROBABILITY
- N independent observations of total pulse eg $\{L, H, L, L \ldots H\}$
- What is probability of observing this?
- We know $P(Low) = \dfrac{exp(\alpha + \beta_s Smoke + \beta_w Weight)}{1 + exp(\alpha + \beta_s Smoke + \beta_w Weight)}$

- So $P(High) = \dfrac{1}{1 + exp(\alpha + \beta_s Smoke + \beta_w Weight)}$

- Also due to independence $P(L, H, H \ldots L)$ can be simplified to $P(L) P(L) P(H) \ldots P(H)$

- Can be expelled a:

$$\prod_{i=1}^{n}\left[\frac{\exp(\alpha + \beta_s\text{Smoke}_i + \beta_w\text{Weight}_i)}{1+\exp(\alpha + \beta_s\text{Smoke}_i + \beta_w\text{Weight}_i)}\right]^{y_i} \times \prod_{i=1}^{n}\left[\frac{1}{1+\exp(\alpha + \beta_s\text{Smoke}_i + \beta_w\text{Weight}_i)}\right]^{1-y_i}$$

- Here $y_i = 1$ if person is $L$ and $y_i = 0$ if person $= H$
- Find vals of $\alpha, \beta_s, \beta_w$ that maximise the probabilities — MLE

## INTREPRETING THE R OUTPUT
- Estimate column gives estimate of $\alpha, \beta_s, \beta_w$    Sign important
- 1 category = low,    ⓪ = high
- Coef of SmokeYes is negative, if individual smoked they would have smaller probability of being assigned to be 1 category → more likely to have H rate
- Coef of Weight is positive, heavier you weigh better chance of low
- Be careful of categorical/binary covariates

Person who Smoked and weighs 190 lb

$$\Rightarrow P(\text{low}) = \frac{\exp(-1.99 - 1.14 \times 1 + 0.03 \times 190)}{1+\exp(-1.99 - 1.19 \times 1 + 0.05 \times 190)} = 0.93$$

If person did not Smoke $\Rightarrow 0.98$

## STANDARD ERRORS
- SE record the uncertainty in the resulting estimate
- Tend to decrease with sample size
- 95% CI    coef estimate $\pm$ 2 (SE of coef)
- $\beta_s \in (-2.30, -0.09)$,    $\beta_w \in (0.001, 0.050)$
- Can be used to find 95% CI for $\exp(\beta_s)$ and $\exp(\beta_w)$
   $\exp(\beta_s) \in (0.1, 0.9)$    $\exp(\beta_w) \in (1.0, 1.05)$

## TESTS
- May wish to test $H_0: \beta_s = 0$  v  $H_1: \beta_s \neq 0$    or $H_0: \beta_w = 0$  v  $H_1: \beta_w \neq 0$
- Compute $\dfrac{(\text{coef est}) - 0}{(\text{SE coef})}$
- Under $H_0$ the quantity is approximately Normal $(0,1)$

- If data had been standardised, then the magnitude of the coef would indicate the importance of the variable in predicting the class
- If we assume true coefficient $\beta$ of a variate was 0, and that actual estimate generated $\hat{\beta}$ was random and subject to a draw from a normal dist with mean 0 and sd to estimated standard error $S_{\hat{\beta}}$ then $\hat{\beta}/S_{\hat{\beta}} \sim N(0,1)$

# MLA

## LOGISTIC REGRESSION
7/04/15

- The column $z$ gives a) $\beta/sd$ and column $Pr(>|z|)$ gives b) the probability of observing a z-value at least that far from 0 under the $N(0,1)$ distributional assumption

- If $Pr(>|z|)$ is small, then under the hyp that true coef is 0 and that the estimate is drawn from a normal distribution with mean 0 and sd as given by S.E., we have observed an unlikely outcome

- If $Pr(>|z|)$ is very small we may start to question the appropriateness of our original hypothesis, a) we are faced with the dilemma of justifying a rare event

- Note that $Pr(>|z|)$ is NOT the probability that the true coef is 0, it is the probability of having observed such a coef estimate under the assumption that the true coef is 0

- For resting pulse data, at 5% significance:
- The estimate for intercept α is not to extreme a) to reject the hyp the true value is 0 - may consider model where it's not included
- Estimate for $\beta_S$ is extreme under H that its true value is 0, may wish to reject such H. If so, we may consider the Smoked indicator a) significant in predicting resting pulse
- Estimate for $\beta_W$ is extreme under the H that its true value is 0, may reject such H. If so, may consider weight value a) significant in predicting resting pulse

## INTERACTIONS
- Including interactions adds possible products of the columns to the model.
- In pulse data: $\text{Log}\left(\frac{p(low)}{1-p(low)}\right) = \alpha + \beta_S \text{Smoked} + \beta_W \text{Weight} + \beta_{SW} \text{Smoked} \times \text{Weight}$

- Interactions allow for effect of one column to be altered depending on value of another column
- In the above, the interaction would be reasonable if the effect that weight has upon resting pulse would differ depending on whether or not the individual smoked.
- For example W of individual may be very important in predicting low pulse if individual smoked but of little importance if individual does not smoke
- When appropriate, interactions can greatly increase model's predictive ability

- Pick model with least AIC — Akaike Information criterion
$$-2\log(\text{likelihood}) + 2p$$

## DEVIANCE

- Consider Saturated model, same number of parameters as data pts
- All 1 or 0.
- Likelihood is $L(Data|1,...,n) = 1$
- Log likelihood is $L(Data|1,...,n) = 0$

- Can be shown that under the $H_0$ hypothesis the proposed model is true, the deviance is approximately distributed as a $\chi^2_{m-n-1}$ distribution, where $n$ is the number of parameters of proposed model

- Hence if deviance is 75% point of the chi-square dist, we may claim under $H_1$ the proposed model is true, we have obtained a rare event

## NULL DEVIANCE

- Simplest model is one in which any data point is assigned a common probability of class assignment regardless of its covariate size
- In this case we have one parameter to determine probability of class assignment and it is estimated as the proportion of data in class "i"
- Deviance for proposed model is called null deviance
- If none of the covariates actually influence class assignment, the model deviance will not be much smaller than null deviance

$$L_{null} \leq L_{mod} \leq L_{max}.$$

- Deviance (residual dev) = $2[Log(L_{max}) - log(L_{mod})]$
- null dev = $2[Log(L_{max}) - log(L_{null})]$

- For n reasonably large and M small, can interpret residual dev as a measure of fit.
- Calculate $1 - pchisq(resid. deviance, deg. freedom)$
- If resulting value is small, then a rare event has occurred under the assumption that our proposed model is correct.

- To test if model is explaining variance in the data we can compute the p-value from a goodness of fit test.

$$1 - pchisq(deviance(res), df. residual(res))$$

- Smaller the value, greater the evidence that there is a significant difference between full model and model of interest. 0.05 cut off assumed
- Often will be a significant difference, so checking the difference of deviance against null deviance against a $\chi^2_p$-dist will indicate whether the model of interest is better than the null model with no covariates

MLA

## LOGISTIC REGRESSION

- Considered as a classification technique when there are two groups
- A new observation is classified as being of one group or another depending on whether the predicted probability falls above or below the threshold of 0.5.
- Similar to LDA.

- In LDA, decision boundary between class $k$ and class $l$ given by

$$Log \frac{P(k|x)}{P(l|x)} = log \frac{\pi_k}{\pi_l} + log \frac{f(x|k)}{f(x|l)} = 0$$

- In LR, by assumption: $Log \frac{P(k|x)}{P(l|x)} = \beta_0 + \beta^T x$

- Have model have same form.
- Difference lie in way linear coefficient are estimated