

# **Section IX**

## **Introduction to Logistic Regression for binary outcomes**

## **Poisson regression**

## Sec 9 - Logistic regression

In linear regression, we studied models where  $Y$  is a continuous variable.

What about the case when  $Y$  is binary?

$Y = 1$  (presence) or  $Y = 0$  (absence)

(Ex:  $Y$  is alive or dead, sick or well, positive or negative)

Linear regression does not work well for two reasons:

1.  $Y$  can't be linearly related to  $X$ 's
2.  $Y$  does not have a Gaussian distribution given  $P$ , the probability that  $Y=1$ , and therefore the errors are not Gaussian.

For (1), we need a linearizing transformation, for (2), we need a different "error" model.

Let  $P$  be the “mean” of the (binary)  $Y$ s for a fixed combinations of  $X$ s.  $P$  is a **proportion**.

$0 \leq P \leq 1$  so  $P$  can't be a linear fcn of  $X$

What about  $P/(1-P) = \text{odds}$

$0 \leq \text{odds} \leq \text{infinity}$  has “floor” of 0

What about  $\log[\text{odds}] = \log[ P/(1-P) ]$

$-\text{infinity} \leq \log[ P/(1-P) ] \leq \text{infinity}$

The  $\log[ P/(1-P) ]$  transformation of  $P$  is called the **logit** (log odds) transformation.

In Logistic regression, we **assume**

$$\text{logit}(P) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

(and that  $Y$  has a bernoulli distribution with mean =  $P$  and  $\text{var}(Y) = P(1-P)$  )

This implies that

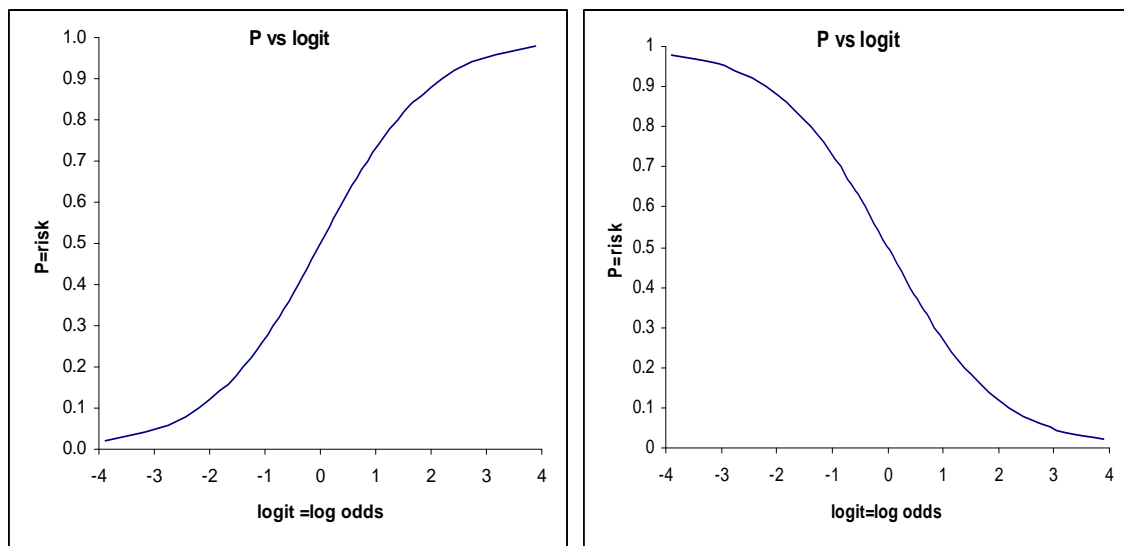
$$[P/(1-P)] = \text{odds} = e^{\text{logit}} = e^{(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}$$

$\Leftrightarrow$

$$\begin{aligned} \text{mean } Y = P &= 1/(1 + e^{-\text{logit}}) \\ &= 1/[1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}] \end{aligned}$$

P is the predicted probability that Y=1 or the predicted proportion of Y's equal to 1 for a given combination of the Xs. P is also the risk, for example, the risk of disease.

Note: If  $\text{odds} = P/(1-P)$  then  $P = \text{odds}/(\text{odds}+1)$



If  $\text{logit} = \log(\text{odds}) =$

$$\ln(P/(1-P)) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

then

$$\text{odds} = (e^{\beta_0}) (e^{\beta_1 X_1}) (e^{\beta_2 X_2}) \dots (e^{\beta_k X_k})$$

or

$$\text{odds} = (\text{base odds}) \text{OR}_1 \text{OR}_2 \dots \text{OR}_k$$

(The base odds is the odds when all  $X$ s equal zero)  
(OR=odds ratio)

Above gives odds for a set of  $X_1, X_2, \dots, X_k$

$$P = \text{risk} = \text{odds}/(\text{odds} + 1) = 1/(1 + \text{odds}^{-1})$$

Or

$$P = 1/(1 + e^{-\log(\text{odds})}) = 1/(1 + e^{-\text{logit}})$$

Or

$$P = 1/[1 + e^{-(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)}]$$

What do the  $\beta$  coefficients mean?

They are the linear rate of change in the logit per unit change in the X. (Huh?)

Example: What if X is coded 0 for males,  
1 for females and

$$\text{logit}(P) = \beta_0 + \beta_1 X$$

If  $X = 0$  (males) then  $\text{logit}(P_m) = \beta_0$

If  $X = 1$  (females) then  $\text{logit}(P_f) = \beta_0 + \beta_1$

therefore  $\beta_1 = \text{logit}(P_f) - \text{logit}(P_m) =$

$$\log\left(\frac{P_f}{1-P_f}\right) - \log\left(\frac{P_m}{1-P_m}\right) =$$

$\log(\text{odds females}) - \log(\text{odds males}) =$

$$\log(\text{odds females/odds males}) =$$

$$\log(\text{Odds ratio}) = \log(\text{OR})$$

$$\text{so, } e^{(\beta_1)} = \text{OR}$$

## Computing Odds

Let  $P$  = proportion with disease

assume  $\text{logit}(P) = \beta_0 + \beta_1 \text{ age} + \beta_2 \text{ sex}$

(Sex = 0 for male, 1 for female)

Odds ratio of disease in females compared to males who are the same age =  $e^{\beta_2}$

Increase in odds of disease for a one year increase in age =  $e^{\beta_1}$  if gender is the same

Increase in odds of disease for 'n' year change in age =  $(e^{\beta_1})^n = e^{n\beta_1}$ .

## Computing odds ratios (ORs) for continuous variables

Outcome: MI= myocardial infarction (coded 0 or 1)

Predictors:

age in years

htn - Hypertension (coded 0 or 1)

smoke – smoking (coded 0 or 1)

$$\text{Logit}(\text{MI}) = \beta_0 + \beta_1 \text{ age} + \beta_2 \text{ htn} + \beta_3 \text{ smoke}$$

What is OR for 40 year old with hypertension vs an otherwise identical 30 year old without hypertension?

$$e^{\beta_0 + \beta_1 40 + \beta_2 (1) + \beta_3 \text{ smoke} - [\beta_0 + \beta_1 30 + \beta_2 (0) + \beta_3 \text{ smoke}]}$$

$$= e^{\beta_1 10 + \beta_2}$$



## Interactions

Let  $Y = \text{CHD}$  (0=no, 1=yes)

$S = \text{smoking}$  (0=no, 1=yes)

$D = \text{drinking alcohol}$  (0=no, 1=yes)

Model:  $\text{logit}(P) = \beta_0 + \beta_1 S + \beta_2 D + \beta_3 S D$

Let **referent category** be:

**no smoking ( $S=0$ ), no drinking ( $D=0$ )**

If  $S=0, D=0$ , odds= $e^{\beta_0}$        $OR_{00} = 1 = e^{\beta_0}/e^{\beta_0}$

If  $S=1, D=0$ , odds= $e^{\beta_0 + \beta_1}$        $OR_{10} = e^{\beta_1}$

If  $S=0, D=1$ , odds= $e^{\beta_0 + \beta_2}$        $OR_{01} = e^{\beta_2}$

If  $S=1, D=1$ , odds= $e^{\beta_0 + \beta_1 + \beta_2 + \beta_3}$        $OR_{11} = e^{(\beta_1 + \beta_2 + \beta_3)}$

When will  $OR_{11} = OR_{10} \times OR_{01}$  ?

If and only if  $\beta_3 = 0$

# Interpretation example

## Potential predictors (13) of in hospital infection mortality (yes or no)

Crabtree, et al JAMA 8 Dec 1999 No 22, 2143-2148

Gender (female or male)

Age - years

APACHE score (0-129)

Diabetes (y/n)

Renal insufficiency / Hemodialysis (y/n)

Intubation / mechanical ventilation (y/n)

Malignancy (y/n)

Steroid therapy (y/n)

Transfusions (y/n)

Organ transplant (y/n)

WBC - count

Max temp - degrees F

Days from admission to treatment

(or > 7 days yes or no)

## Interpretation example

**Table 3. Stepwise Logistic Regression Factors Associated With Mortality for All Infections**

<b>Characteristic</b>	<b>Odds Ratio (95% CI)</b>	<b>p value</b>
Incr APACHE score	1.15 (1.11-1.18)	<.001
Transfusion (y/n)	4.15 (2.46-6.99)	<.001
Increasing age	1.03 (1.02-1.05)	<.001
Malignancy	2.60 (1.62-4.17)	<.001
Max Temperature	0.70 (0.58-0.85)	<.001
Adm to treat>7 d	1.66 (1.05-2.61)	0.03
Female (y/n)	1.32 (0.90-1.94)	0.16

\*APACHE = Acute Physiology & Chronic Health Evaluation Score

## Diabetes complications -Descriptive stats

"Obese"=actual weight above expected weight

Table of obese by diabetes complication

obese            diabetes complication

Freq	no- 0	yes- 1	Total	% yes
no 0	56	28	84	28/84=33%
yes 1	20	41	61	41/61=67%
Total	76	69	145	

%obese        26%        59%        (RR=2.0, OR=4.1)  
 P value < 0.001

Fasting glucose ("fastglu") mg/dl

	n	min	median	mean	max
No complication	76	70.0	90.0	91.2	112.0
Complication	69	75.0	114.0	155.9	353.0

P value =

Steady state glucose ("steadyglu") mg/dl

	n	min	median	mean	max
No complication	76	29.0	105.0	114.0	273.0
Complication	69	60.0	257.0	261.5	480.0

P value =



## Diabetes complications-Logistic model (SAS output)

Parameter	DF	beta	SE(b)	Chi-Square	p
Intercept	1	-14.702	3.231	20.706	<.0001
obese	1	0.328	0.615	0.285	0.5938
fastglu	1	0.108	0.031	12.456	0.0004
steadyglu	1	0.0225	0.0053	18.322	<.0001

### Odds Ratio Estimates

Effect	Point Estimate	95% Wald Confidence Limits
obese	1.388	0.416 4.631
fastglu	1.114	1.049 1.182
steadyglu	1.023	1.012 1.033

Log odds diabetes complication =

$$-14.7 + 0.328 \text{ obese} + 0.108 \text{ fast glu} + 0.023 \text{ steady glu}$$

Q-What is the model predicted “risk” if you are obese, have fasting glucose=100 mg/dl and a steady state glucose=120 mg/dl ?

$$A\text{-logit} = -14.7 + 0.328 + 0.108(100) + 0.023(120) = -0.872$$

$$\text{Odds} = \exp(-0.872) = 0.418, \quad P = 0.418 / 1.418 = 0.295$$

## Statistical Significance of the $\beta$ s

Just as there are t statistics for the  $\beta$  coefficients in linear regression, the Wald chi-square statistic is computed for each  $\beta$  in logistic regression. The Wald chi-square is defined as

$$\chi^2_{\text{wald}} = [\beta / \text{SE}(\beta)]^2$$

A confidence interval for the true  $\beta$  is formed in the usual way

$$\text{CI: } \beta \pm Z \text{ SE}(\beta)$$

where Z is the appropriate Gaussian percentile. For 95% intervals,  $Z=1.96$  as usual.

First form (95%) CI for  $\beta$  on log scale

$$b - 1.96 \text{ SE}, \quad b + 1.96 \text{ SE}$$

then take antilogs of each end to get confidence interval for the  $\text{OR} = e^b$

$$e^{[b - 1.96 \text{ SE}]}, \quad e^{[b + 1.96 \text{ SE}]}$$

## Model fit-Linear vs Logistic regression

k variables, n observations

<u>Variation</u>	<u>df</u>	<u>sum square or deviance</u>
Model	k	G
Error	n-k	D
Total	n-1	T <-fixed

$Y_i = i^{\text{th}}$  observation,  $\hat{Y}_i = \text{prediction for } i^{\text{th}} \text{ obs}$

<u>statistic</u>	<u>Linear regr</u>	<u>Logistic regr</u>
$D/(n-k)$	Residual $SD_e$	Mean deviance
$\Sigma[(Y_i - \hat{Y}_i)/\hat{Y}]^2$	--	Hosmer-L $\chi^2$
$\text{Corr}(Y, \hat{Y})^2$	$R^2$	Cox-Snell $R^2$
G/T	$R^2$	Pseudo $R^2$

Good regression models have large G and small D. For logistic regression,  $D/(n-k)$ , the mean deviance, should be near 1.0.



## Model fit-Analysis of deviance

### Diabetes complication model

	<u>df</u>	<u>-2log L</u>	<u>p value</u>
Model (G)	3	117.21	< 0.001
<u>Error (deviance=D)</u>	<u>141</u>	<u>83.46</u>	<u>0.99</u>
Total (G+D)	144	200.67	

$$\text{Mean deviance} = 83.46/141 = 0.59$$

Want mean deviance  $\leq 1.0$  if model fits the data.

$$R^2_{\text{pseudo}} = 117.21/200.67 = 0.584$$

$$R^2_{\text{cox-snell}} = 0.554$$

## Model fit–Diabetes complication data

### Hosmer-Lemeshow $\chi^2$ Goodness of fit test

$$\chi^2 = (\text{obs-expected})^2/\text{expected}$$

The observed and model based expected frequency values (ie. number of persons) by deciles

Group	Total	complication = 1		complication = 0	
		Observed	Expected	Observed	Expected
1	16	0	0.23	16	15.77
2	15	0	0.61	15	14.39
3	15	0	1.31	15	13.69
4	15	4	2.48	11	12.52
5	15	8	4.53	7	10.47
6	15	6	8.49	9	6.51
7	15	13	12.74	2	2.26
8	16	15	15.60	1	0.40
9	23	23	23.00	0	0.00
<b>total</b>	<b>145</b>	<b>69</b>	<b>69.00</b>	<b>76</b>	<b>76.00</b>

Chi-Square	DF	p value
9.8949	7	0.1946

If the model fits the data (the null hypothesis here), the H-L chi-square statistic should be small and the p value should be large.

## Goodness of fit versus variation accounted for

It is possible that that  $D/(n-k)$ , the mean deviance, and the Hosmer-Lemeshow goodness of fit  $\chi^2$  (H-L  $\chi^2$ ) can both indicate that the model “fits” the data (mean deviance near 1.0, H-L  $\chi^2$  small and its p value large) and yet the corresponding  $R^2$  values for the same model may be low, implying little of the variation is accounted for. Is this a contradiction?

No. The mean deviance and H-L  $\chi^2$  are based on comparing observed versus expected frequencies across the “J” covariate patterns (“cells”) where these J covariate patterns are **made using the variables in the model**. This is like comparing observed to expected cell MEANS in linear regression.

In contrast, the  $R^2$  statistics are based on how much of the total variation is accounted for by the model, not just how well the model fits the means.

In effect, the  $R^2$ , as in linear regression, is not just assessing fit across the current J covariate patterns, but is assessing how much the total variation of Y (including variation within the J covariate cells) is accounted for by the current model. This variation may be greater than the variation of the (average) Y across the J covariate patterns.

So, if the mean deviance and H-L  $\chi^2$  are satisfactory but the  $R^2$  is low, **adding interactions or other terms to the model using the current variables will not improve the  $R^2$**  much. This result implies that **there are other factors not currently in the model**, that are needed to account for the variation of Y.

## Goodness of fit–Logistic & linear regression

### Logistic (hosmer-lemeshow)

$P$ vs $\hat{P}$	$(1-P)$ vs $(1-\hat{P})$

### Linear

$Y$ vs $\hat{Y}$	

For a given covariate pattern, the model predicted values and the mean values can agree (good fit) even though the predicted value and the individual values do not agree as well. That is, the predicted values can agree with the observed values “on average” but not for each individual.

## Sensitivity & Specificity for Logistic

If one has a binary outcome like positive or negative (say, for a disease), one can imagine classifying persons into these two groups on the basis of their X values and then comparing the classification to their actual status.

The sensitivity and specificity are measures of how good the classification is:

	True pos	True neg
Classify pos	a	b
Classify neg	c	d
total	a+c	b+d

Sensitivity =  $a/(a+c)$ , specificity =  $d/(b+d)$

false neg =  $c/(a+c)$ , false pos =  $b/(b+d)$

false neg =  $1 - \text{sensitivity}$

false pos =  $1 - \text{specificity}$

“accuracy =  $W \text{ sensitivity} + (1-W) \text{ specificity}$ ”

A good classification rule (& thus a good logistic model) has high sensitivity and high specificity.

If we choose a cutpoint,  $P_c$ , we can use the predicted values ( $P_i$ ) from our logistic rule to classify persons.

Classify as positive if  $P_i = \text{predicted } P > P_c$

Otherwise classify as negative.

For a given  $P_c$  we can compute sensitivity and specificity and also an overall penalty function, of the form

$$\text{penalty} = W(\text{false neg}) + (1-W)(\text{false pos})$$

where  $W$  gives the relative weights we assign to a false positive or a false negative.

We can choose  $P_0$  to minimize the value of the penalty function. This is the same as maximizing

$$\text{accuracy} = W(\text{sensitivity}) + (1-W)\text{specificity}$$

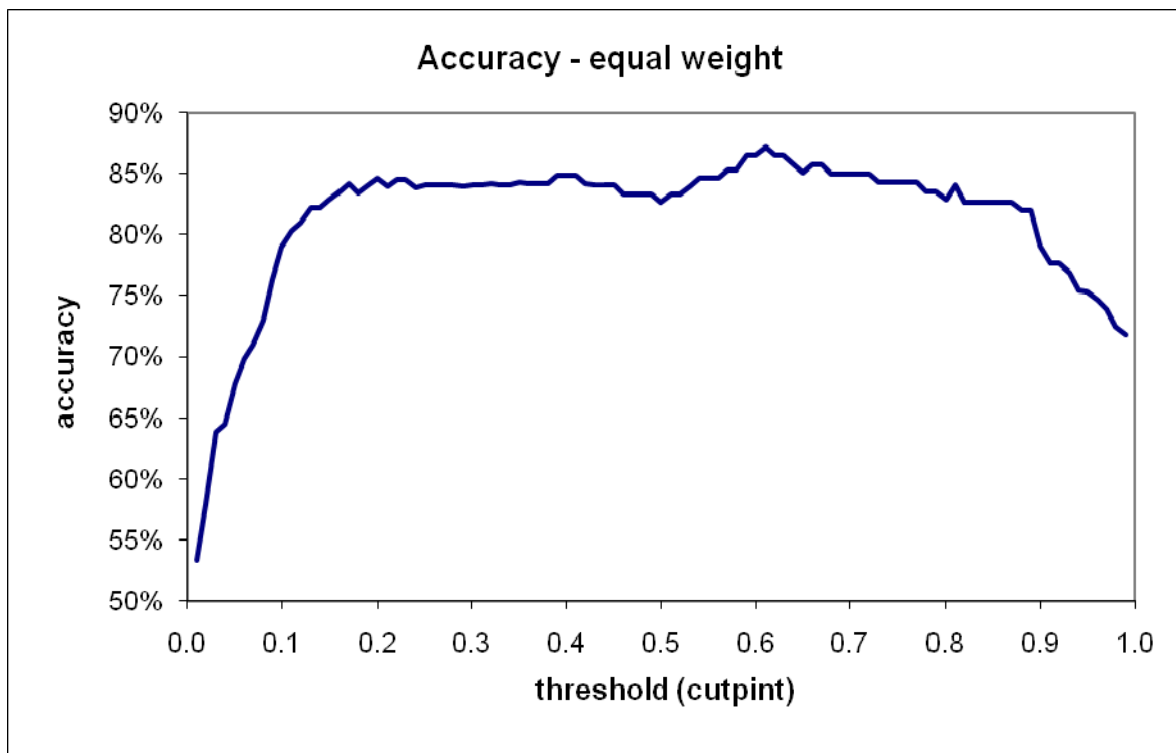
$W=0.5$  implies equal weight.

## Accuracy: diabetes complication model

$\text{logit}(P_i) = -14.7 + 0.328 \text{ obese} + 0.108 \text{ fast glu} + 0.023 \text{ steady glu}$

$$P_i = e^{\text{logit}} / (1 + e^{\text{logit}})$$

Compute  $\text{logit}(P_i)$  (or  $P_i$ ) for all observations and do an ROC as below. Find value of  $\text{logit}$  or  $P_i$  (called  $P_0$ ) that maximizes  $\text{accuracy} = (\text{sensitivity} + \text{specificity}) / 2$



### Best accuracy

$$\text{Logit} = 0.447, P_0 = (e^{0.447}) / (1 + e^{0.447}) = 0.61$$

	True comp	True no comp
Predicted yes	<b>55</b>	4
Predicted no	14	<b>72</b>
Total	69	76

$$\text{Sensitivity} = 55/69 = 79.7\%$$

$$\text{Specificity} = 72/76 = 94.7\%$$

overall accuracy (equal weight)

$$= 0.5 (.797) + 0.5(.947) = 87.2\%$$

However, it is possible that a model can fit the data by the H-L goodness of fit  $\chi^2$  but still have poor sensitivity and specificity. This implies that most subjects are at intermediate risk, not high and/or low risk.



## **C statistic – summary accuracy**

The C (concordance) statistic is the **area under the ROC curve** and varies from 0.5 (worst) to 1.0 (best).

Imagine forming all  $n_0 \times n_1$  possible pairs of the  $n_0$  obs with  $Y=0$  and the  $n_1$  obs with  $Y=1$ .

If the logit for the pair member with the 1 is higher than the logit for the pair member with the 0, call this pair a “concordant” pair. Otherwise, call it discordant unless their logits are identical, in which case, call it a “tie”.

$$C = \frac{\text{num concordant} + 0.5 \text{ num ties}}{n_0 n_1}$$

If  $Y=1$  is disease and  $Y=0$  is non disease, C is the probability that, in any pair, the diseased subject has a higher predicted probability of disease than the non diseased subject.

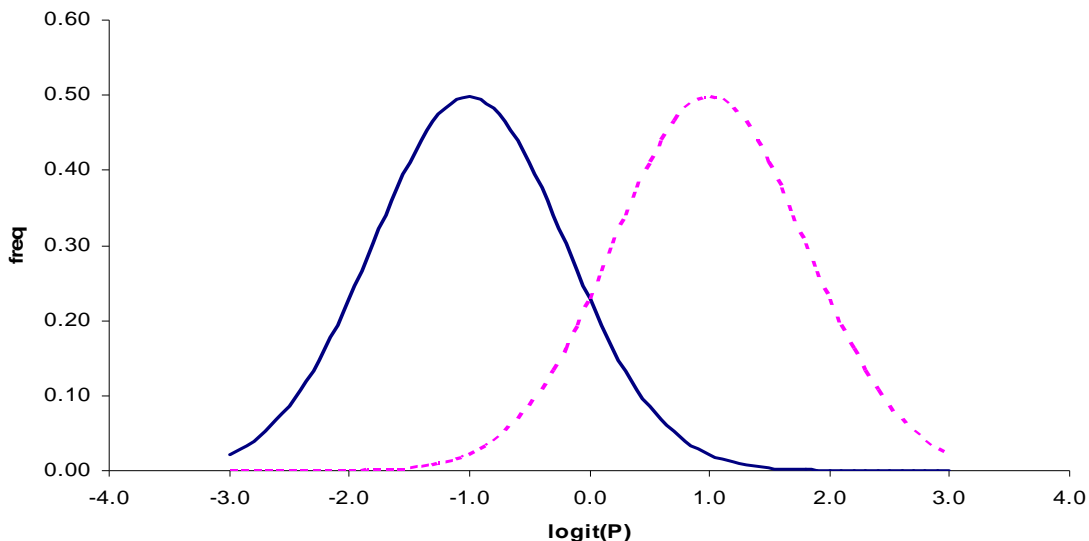
For Diabetes comp,  $C = 0.949$

## Logistic model is a discrimination model

The logit score from logistic regression can also be considered a discrimination “score”.

One can plot histograms of the logit score in those who have  $Y=1$  versus those who have  $Y=0$ . The best cutpoint or threshold that separates these two histograms corresponds to  $\text{logit}(P_0)$ , the logit of  $P_0$ .

Examining the histograms of the logits is a generalization of comparing the distribution of any continuous variable  $X$  between the  $Y=1$  versus  $Y=0$  groups.



# Poisson regression

When  $Y$  is a positive integer (ie  $Y=0,1,2,3\dots$ ), we cannot model  $Y$  or its mean as a linear function of  $X$ s since  $Y$  has a “floor” of zero.

$$0 \leq Y \leq \text{infinity}$$

We model  $\log(\hat{Y}) = \log(\mu)$  as a linear function of the  $X$  variables so that

$$\text{Mean } Y = \hat{Y} = \mathbf{exp}(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)$$

In this way, Mean  $Y = \hat{Y}$  and  $Y$  can never be less than zero. This model is the **Poisson Regression** model.

## Poisson Regression interpretation of $\beta$ s

For poisson regression,

$$\ln(\hat{Y}) = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k$$

so,  $\beta_i$  is the rate of change of  $\ln(\hat{Y})$  per unit change in  $X_i$  holding the other  $X$ s constant.

$$\text{Since } \hat{Y} = \mathbf{exp}(\beta_0 + \beta_1 X_1 + \dots + \beta_k X_k)$$

$$d\hat{Y}/dX_i = \beta_i \hat{Y}, \quad \text{or} \quad \beta_i = [d\hat{Y}/dX_i] / \hat{Y}$$

so,  $100 \beta_i$  is the percent change in  $\hat{Y}$  per unit change in  $X_i$ , holding all else constant. For  $X_i$ ,  $\exp(\beta_i)$  is the **mean ratio** of  $\hat{Y}$  for  $X_i$  to  $\hat{Y}$  for  $X_i-1$ . This is similar to the **odds ratio** in logistic regression.

# Linear vs Logistic regression

Method	Linear regression	Logistic regression
Outcome (Y)	Continuous / means	Binary / Proportions
Predictors (Xs)	Continuous and discrete*	Continuous and discrete
Overall significance test	F test= MS model / MS error	Chi-square test= Reduced-Full
Model Fit stats	R <sup>2</sup> , residual SD	Mean deviance=1, fit test, Hosmer-Lemeshow stat
Model variation	Model Sum of Squares	“Difference” log likelihood
Error variation	Residual Sum of Squares	“Full” log likelihood= <b>deviance</b>
Total variation= Model + Error	Total SS = Model SS + Residual SS	“Reduced = Difference + Full

\* The analysis of variance (ANOVA) model is a special case of linear regression where all the predictors are discrete.

## Simplified Regression guide

$$\text{let } \mathbf{XB} = b_0 + b_1X_1 + b_2X_2 + \dots + b_kX_k$$

outcome      range      method      link      error

continuous      -infinity to + infinity      linear regr /ANOVA       $y=\mathbf{XB}$       gaussian

(if all of the predictors are discrete, do ANOVA  
if there is a random person effects, do repeated measure ANOVA)

binary      0 to 1      logistic regr.       $y=1/(1+\exp(-\mathbf{XB}))$       binomial

positive integer      0 to infinity      poisson regr,       $y=\exp(\mathbf{XB})$       poisson

If the outcome is continuous but non linear in the beta parameters, must do non linear regression with Gaussian errors.

## Analysis of deviance (log likelihood)

### Depression model

	<u>df</u>	<u>-2log L</u>	<u>p value</u>
Model (G)	3	16.14	< 0.001
<u>Error (deviance=D)</u>	<u>290</u>	<u>252.00</u>	
Total (G+D)	293	268.14	

$$\text{Mean deviance} = 252/290 = 0.869$$

Want mean deviance < 1.0

$$\text{Pseudo } R^2 = 16.14/268.14 = 0.060$$

$$R^2 \text{ cox-snell} = 0.0893$$

Hosmer Lemeshow chi-square =