

SAS - Logistic - exs. pdf

15	16	17	18	19	20	21	22	23	24	25
12.0833	0.16009	0.1806104	-0.53576	-0.544155	0.924241	0.957464	-0.924241	-0.957464	-0.924241	-0.957464
12.0833	0.93197	0.948491	-1.40335	-1.367538	0.248491	0.248491	-1.40335	-1.367538	0.248491	0.248491
15.0833	0.35902	0.9671908	-0.506855	-0.488336	0.9671908	0.9671908	-0.506855	-0.488336	0.9671908	0.9671908
15.0833	0.39000	0.3326479	1.2172483	1.2016944	0.3326479	0.3326479	1.2172483	1.2016944	0.3326479	0.3326479
15.0833	0.39336	0.973399	-0.1025	-0.095574	0.973399	0.973399	-0.1025	-0.095574	0.973399	0.973399
15.0833	0.47242	0.4283346	0.93569	0.918236	0.47242	0.47242	0.93569	0.918236	0.47242	0.47242
15.0833	0.37247	0.828602	-0.828602	-0.799296	0.37247	0.37247	-0.828602	-0.799296	0.37247	0.37247
15.0833	0.38246	0.9361234	-0.828602	-0.799296	0.38246	0.38246	-0.828602	-0.799296	0.38246	0.38246
17.58	1.00000	0.3994267	0.7756656	1.0968278	0.7756656	0.7756656	0.7756656	1.0968278	0.7756656	0.7756656

Leukemia Data

Feigl and Zelen reported the survival time in weeks and the white cell blood count (WBC) at time of diagnosis for 33 patients who eventually died of acute leukemia. Each person was classified as AG+ or AG-, indicating the presence or absence of a certain morphological characteristic in the white cells. Four variables are given in the SAS data set: WBC, a binary indicator variable IAG (1 for AG+, 0 for AG-), NTOTAL (the # of patients with the given combination of IAG and WBC), and NRES (the # of NTOTAL that survived at least one year from the time of diagnosis).

The researchers are interested in modeling the probability p of surviving at least one year as a function of WBC and IAG. They believe that WBC should be transformed to a log scale, given the skewness in the WBC values.

As an initial step in the analysis, consider the following model (with multiple effects):

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \text{LWBC} + \beta_2 \text{IAG},$$

LWBC = log of WBC
IAG = dummy

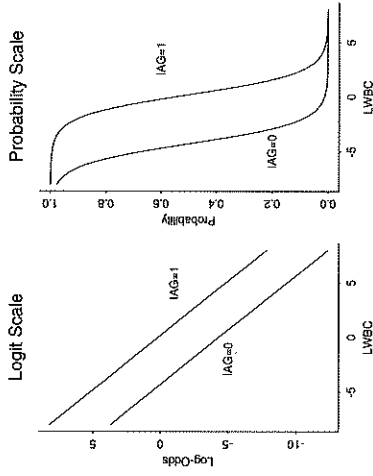
where LWBC = log WBC. The model is best understood by separating the AG+ and AG- cases. For AG- individuals, IAG=0 so the model reduces to

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \text{LWBC} + \beta_2 * 0 = \beta_0 + \beta_1 \text{LWBC}.$$

For AG+ individuals, IAG=1 and the model implies

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \text{LWBC} + \beta_2 * 1 = (\beta_0 + \beta_2) + \beta_1 \text{LWBC}.$$

The model without IAG (i.e. $\beta_2 = 0$) is a simple logistic model where the log-odds of surviving one year is linearly related to LWBC, and is independent of AG. The reduced



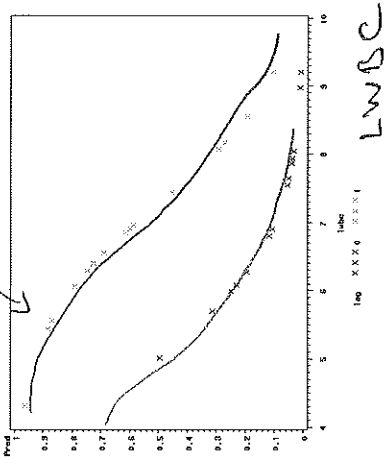
model with $\beta_2 = 0$ implies that there is no effect of the AG level on the survival probability once LWBC has been taken into account.

Including the binary predictor IAG in the model implies that there is a linear relationship between the log-odds of surviving one year and LWBC, with a constant slope for the two AG levels. This model includes an effect for the AG morphological factor, but more general models are possible. A natural extension would be to include a product or interaction effect, a point that I will return to momentarily.

The parameters are easily interpreted: β_0 and $\beta_0 + \beta_2$ are intercepts for the population logistic regression lines for AG- and AG+, respectively. The lines have a common slope, β_1 . The β_2 coefficient for the IAG indicator is the difference between intercepts for the AG+ and AG- regression lines. A picture of the assumed relationship is given above for $\beta_1 < 0$. The population regression lines are parallel on the logit scale only, but the order between IAG groups is preserved on the probability scale.

Before looking at SAS output for the equal slopes model, note that the data set has 30 distinct IAG and LWBC combinations, or 30 "groups" or samples. Only two samples have more than 1 observation. The majority of the observed proportions surviving at least one year (number surviving ≥ 1 year / group sample size) are 0 (i.e. 0/1) or 1 (i.e. 1/1).

$IAG = \begin{cases} 1 & \text{red (AG+)} \\ 0 & \text{black (AG-)} \end{cases}$



predicted probability (survival) \hat{p}

The ratios of D and X^2 divided by the df are both less than 1, indicating that there are no gross deficiencies with the model. Recall that the **Analysis Of Parameter Estimates** table gives p-values for testing the hypothesis that the regression coefficients are zero for each predictor in the model. The two predictors are LWBC and IAG, so the small p-values indicate that LWBC and IAG are important predictors of survival in this model. The LR test p-values in the LR table lead to the same conclusion. If either predictor was insignificant, I would consider refitting the model omitting the least significant effect, as in regression.

Given that the model fits reasonably well, a test of $H_0: \beta_2 = 0$ might be a primary interest here. This checks whether the regression lines are identical for the two AG levels, which is a test for whether AG affects the survival probability, after taking LWBC into account. This test is rejected at any of the usual significance levels, suggesting that the AG level affects the survival probability (assuming a very specific model).

A plot of the predicted survival probabilities as a function of LWBC is given with the SAS output, using different colors to identify the two IAG groups. The model indicates that the probability of surviving at least one year from the time of diagnosis is a decreasing function of LWBC. For a given LWBC the survival probability is greater for AG+ patients

than for AG- patients. This tendency is consistent with the observed proportions, which show little information about the exact form of the trend.

The estimated (ML) survival probabilities satisfy

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = 5.54 - 1.11LWBC + 2.52IAG.$$

For AG- individuals with IAG=0, this reduces to

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = 5.54 - 1.11LWBC,$$

or equivalently,

$$\hat{p} = \frac{\exp(5.54 - 1.11LWBC)}{1 + \exp(5.54 - 1.11LWBC)}.$$

For AG+ individuals with IAG=1,

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = 5.54 - 1.11LWBC + 2.52 * (1) = 8.06 - 1.11LWBC,$$

or

$$\hat{p} = \frac{\exp(8.06 - 1.11LWBC)}{1 + \exp(8.06 - 1.11LWBC)}.$$

The equivalence between logit and probability scales is that $\log(p/(1-p)) = z$ implies $p = \exp(z)/(1 + \exp(z))$, regardless of how many effects define z .

Using the **logit scale**, the difference between AG+ and AG- individuals in the estimated log-odds of surviving at least one year, at a fixed but arbitrary LWBC, is the estimated IAG regression coefficient:

$$AG^+ - AG^- = (8.06 - 1.11LWBC) - (5.54 - 1.11LWBC) = 2.52 = \log\left(\frac{\pi^+}{\pi^-}\right)$$

Using properties of exponential functions, the odds that an AG+ patient lives at least one year is $\exp(2.52) = 12.42$ times larger than the odds that an AG- patient lives at least one year, regardless of LWBC. This statistical summary is called the **adjusted odds ratio** for the IAG effect. A CI for the adjusted odds ratio is obtained by exponentiating the CI for the corresponding regression effect given in the parameter estimates table: $\exp(.38) = 1.47$ to $\exp(4.66) = 105.35$.

$$2.52 = \log(\pi^+) - \log(\pi^-) \Rightarrow \pi^+ = \frac{p^+}{1-p^+} \approx \frac{0.95}{0.05} = 19$$

$$= \log\left(\frac{\pi^+}{\pi^-}\right) \Rightarrow \pi^- = \frac{p^-}{1-p^-} \approx \frac{0.05}{0.95} = 0.0526$$

$$\Rightarrow e^{2.52} = 12.42 = \frac{\pi^+}{\pi^-} \Rightarrow \pi^+ = 12.42 \pi^-$$

$$\Rightarrow 19 = 12.42 \pi^- \Rightarrow \pi^- = 1.53\%$$

$$\Rightarrow \pi^+ = 19 \times 1.53\% = 29.07\%$$

odds S decreases by 67% for each unit increase in LWBC.

$\pi^{+1} = \frac{1}{3} \pi^0$

Similarly, the estimated odds ratio of .33 = $\exp(-1.11)$ for LWBC implies that the odds of surviving at least one year is reduced by a factor of 3 for each unit increase of LWBC.

Although the equal slopes model appears to fit well, a more general model might fit better. A natural generalization here would be to add an interaction, or product term, IAG * LWBC to the model. The logistic model with an IAG effect and the IAG * LWBC interaction is equivalent to fitting separate logistic regression lines to the two AG groups. This interaction model provides an easy way to test whether the slopes are equal across AG levels. I will note that the interaction term is not needed here.

Data from a Budworm Experiment

The purpose of this experiment was to assess the resistance of tobacco budworm larvae to cypermethrin. Batches of 20 pyrethroid-resistant moths (i.e. moths that have grown resistant to the pesticide pyrethroid) of each sex were exposed to a range of doses of cypermethrin two days after emergence from pupation. The number of moths which were either knocked down (uncoordinated movements) or dead was recorded 72 hours after treatment. We will model the probability of uncoordinated movement, or death, as a function of sex and dose.

The SAS program below inputs the number of responders for each combination of sex and dose. The sample size of 20 for each sex-dose combination is defined directly in the data step, rather than entered with the data. A plot of the observed proportion of responders against the log(dose) (log-dose is traditionally used in bioassay settings such as this) shows that increased doses are more effective, and that males tend to be more affected than females.

The first `genmod` proc fits a logistic model with main effects for sex (a factor) and log(dose) (a predictor) and a sex-by-log(dose) interaction. The `model` statement specifies the response in the number of events divided by sample size format as used in the `menarche` data. The `class` statement identifies sex as a factor.

You should probably read the ANCOVA (analysis of covariance) notes from my Data Analysis II website if you are not familiar with how factors are treated in regression models. Factors are typically used with qualitative effects such as sex, ethnicity and so on. The (main) effect for a factor with c levels (sex has $c = 2$ levels) is captured through c indicator (binary,

or 0-1) variables, one for each level of the factor. However, perfect collinearity results from including all c binary predictors in a model with an intercept. Without loss of generality, one of the binary predictors can be removed from the model. The group corresponding to the omitted binary predictor is called the **baseline** group. Interactions between a factor and a predictor (or another factor) corresponds to adding one df effect for each product of the binary variables with the other feature in the interaction. If a factor has $c > 2$ levels the SAS output will include estimates of the regression coefficients for the $c - 1$ single df effects while the LR table will include a test of significance for the entire effect (i.e. a single $c - 1$ df test on the importance of the effect).

By default, SAS uses the category of the class variable with the highest alphanumeric value as the baseline group. Thus, in the budworm analysis, males (M) are the baseline group, and the sex effect corresponds to a single binary predictor with levels 1 for females and 0 for males. The interaction between ld and sex corresponds to a predictor that is the product of ld and the binary sex effect. (ld = log(dose))

Letting p_j be the success probability (i.e. probability of death or uncoordinated movement) for batch j , and defining $sex_j = 1$ if batch j consists of females and 0 if batch j consists of males (baseline), the fitted model is equivalent to:

$$\log\left(\frac{p_j}{1-p_j}\right) = \beta_0 + \beta_1 sex_j + \beta_2 ld_j + \beta_3 sex_j * ld_j$$

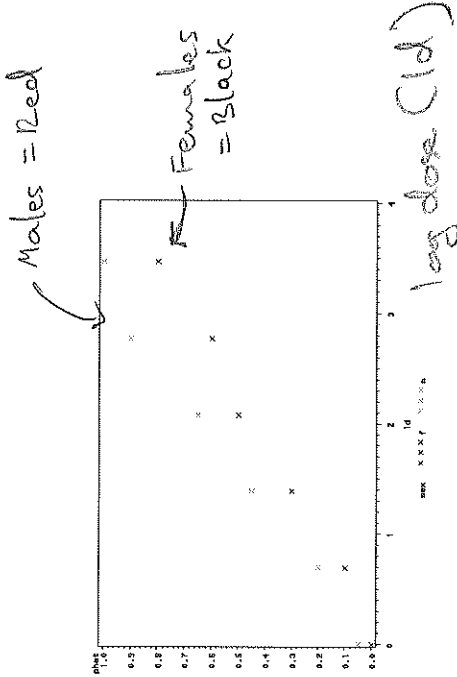
where ld_j is the log-dose given to batch j . The model implies that the log-odds of a response is linearly related to log(dose), but that males and females have distinct lines. To see this, note that the model reduces to

$$\log\left(\frac{p_j}{1-p_j}\right) = \beta_0 + \beta_2 ld_j \quad \text{sex} = 0$$

for males (sex = 0) and

$$\log\left(\frac{p_j}{1-p_j}\right) = \beta_0 + \beta_1 + (\beta_2 + \beta_3) ld_j \quad \text{sex} = 1$$

for females. Thus, β_0 and β_2 give the intercept and slope for the baseline group, males, and β_1 and β_3 give the difference in intercepts and slopes between females and males.



```

proc gplot;
  plot phat*ld=sex;
run;

Model with unequal slopes
-----
proc genmod data=d1;
  classes sex;
  model nresp/ntot = sex ld sex*ld/dist=bin
        link=logit noscale type3;
run;

```

Model Information

Data Set WORK.D1
 Distribution Binomial
 Link Function Logit
 Response Variable (Events) nresp
 Response Variable (Trials) ntot

Number of Observations Read 12
 Number of Observations Used 11
 Number of Events 11
 Number of Trials 240

Class Level Information

Class	Levels	Values
sex	2	f m

Criteria For Assessing Goodness Of Fit

Criterion	DF	Value	Value/DF
Deviance	8	4.9837	0.6242
Scaled Deviance	8	4.9837	0.6242
Pearson Chi-Square	8	3.5047	0.4381
Scaled Pearson X2	8	3.5047	0.4381
Log Likelihood		-105.7388	

Algorithm converged.

Analysis Of Parameter Estimates

Parameter	DF	Estimate	Standard Error	Wald 95% Confidence Limits	Chi-Square	Pr > ChiSq
Intercept	1	-2.8186	0.5480	-3.8926 -1.7445	26.46	<.0001
sex	1	-0.1750	0.7783	-1.7004 1.3505	0.05	0.8221
sex	0	0.0000	0.0000	0.0000 0.0000		
ld	1	1.8163	0.3059	1.2166 2.4159	35.24	<.0001

It is important to note that the choice of a baseline group for a factor impacts the interpretation of the parameters, but does not impact the interpretation of the model. Provided the model follows the hierarchy principal. For example, how would the parameters and model change if we ran the same program but had defined sex to have 2 levels: M for males and W for females?

We could have treated IAG as a factor in the leukemia analysis (even though it was coded numerically), but SAS would have treated IAG=1 as the baseline group. Treating IAG as a predictor (not a class variable) essentially leads to IAG=0 as the baseline group or level. You should read the online help for the class statement in **genmod** to learn how to change the default choice for the baseline group.

```

options ls=79 nodate;
data d1;
  input sex $ dose nresp @@;
  ntot = 20;
  phat = nresp/ntot;
  ld = log(dose);
cards;
  1 m 2 4 m 4 9 m 8 13
  1 f 16 2 f 2 f 4 6 f 8 10
  ;
symbol1 value=x interpol=none width=2;
symbol2 value=y interpol=join width=2;

```

```

ld*sex      f      1      -0.5091      0.3895      -1.2726      0.2643      1.71      0.1912
ld*sex      m      0      0.5000      0.6000      0.6000      0.6000
Scale       0      1.0000      0.0000      1.0000      1.0000
NOTE: The scale parameter was held fixed.

```

LR Statistics For Type 3 Analysis

```

Source      DF      Chi-Square      Pr > ChiSq
sex         1         0.95          0.8221
ld         1        112.73         <.0001
ld*sex      1         1.76          0.1842

```

Looking at the output, we see

1. from the data plot that the pesticide is more effective as the dose increases
2. the effect of dose is statistically significant ($p < .0001$)
3. the sex-by-log(dose) interaction is not significant ($p = .18$ on LR test) so an equal slopes model (i.e. the original model with no interaction) might be considered
4. the sex effect (which estimates the difference in intercepts) is not significant in the presence of the interaction (but is significant if interaction is excluded)
5. The Deviance and Pearson statistics are smaller than the degrees of freedom, which indicate that the model fits the data. A chi-squared approximation is appropriate here, and the p-values would be large.

The fitted model is

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = -2.82 - .18sex_j + 1.82ld_j - .51sex_j * ld_j$$

which reduces to

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = -2.82 + 1.82ld_j$$

for males and

$$\log\left(\frac{\hat{p}_i}{1-\hat{p}_i}\right) = -2.82 - .18 + (1.82 - .51)ld_j = -3.00 + 1.31ld_j$$

for females.

Recalling that $\log(p/(1-p)) = z$ implies $p = \exp(z)/(1 + \exp(z))$, we can write the fitted model as

$$\hat{p}_j = \frac{\exp(-2.82 + 1.82ld_j)}{1 + \exp(-2.82 + 1.82ld_j)} \quad \text{M}$$

for males and

$$\hat{p}_j = \frac{\exp(-3.00 + 1.31ld_j)}{1 + \exp(-3.00 + 1.31ld_j)} \quad \text{F}$$

for females.

A reasonable next step would be to fit a model with equal slopes (i.e. no sex-by-ld interaction), which is done below. I will leave to you to decide what would be the best final model and how to summarize the analysis.

Equal slopes model

```

proc genmod data=d1;
class sex;
model nresp/ntot = sex ld/dist=bin link=logit noscale type3;
run;

```

Criteria For Assessing Goodness Of Fit

Criterion	DF	Value	Value/DF
Deviance	9	6.7571	0.7508
Pearson Chi-Square	9	5.3060	0.5896

Analysis Of Parameter Estimates

Parameter	DF	Estimate	Standard Error	Wald 95% Confidence Limits	Chi-Square	Pr > ChiSq
Intercept	1	-2.3724	0.3855	-3.1280 -1.6168	37.87	<.0001

Using PROC LOGISTIC

This handout illustrates the logistic procedure in SAS - the document is old, and has not been updated with nicer graphics and newer features in SAS.

The UNM Trauma Data (Model survival probability)

The data to be analyzed here were collected on 3132 patients admitted to The University of New Mexico Trauma Center between the years 1991 and 1994. For each patient, the attending physician recorded their age, their revised trauma score (RTS), their injury severity score (ISS), whether their injuries were blunt (i.e. the result of a car crash: BP=0) or penetrating (i.e. gunshot wounds: BP=1), and whether they eventually survived their injuries (SURV=0 if not, SURV=1 if survived). Approximately 10% of patients admitted to the UNM Trauma Center eventually die from their injuries.

The ISS is an overall index of a patient's injuries, based on the approximately 1300 injuries cataloged in the Abbreviated Injury Scale. The ISS can take on values from 0 for a patient with no injuries to 75 for a patient with 3 or more life threatening injuries. The ISS is the standard injury index used by trauma centers throughout the U.S. The RTS is an index of physiologic injury, and is constructed as a weighted average of an incoming patient's systolic blood pressure, respiratory rate, and Glasgow Coma Scale. The RTS can take on values from 0 for a patient with no vital signs to 7.84 for a patient with normal vital signs.

Champion et al. (1981) proposed a logistic regression model to estimate the probability of a patient's survival as a function of RTS, the injury severity score ISS, and the patient's age, which is used as a surrogate for physiologic reserve. Subsequent survival models included the binary effect BP as a means to differentiate between blunt and penetrating injuries.

We will develop a logistic model for predicting survival from ISS, AGE, BP, and RTS. Data on the number of severe injuries in each of the nine body regions is also included in the database, so we will also assess whether these features have any predictive power. The following labels were used to identify the number of severe injuries in the nine regions: AS

sex	f	1	-1.1007	0.3558	-1.7982	-0.4033	9.57	0.0020
sex	m	0	0.0000	0.0000	0.0000	0.0000		
ld	1	1.5353	0.1891	1.1647	1.9060	65.92		<.0001

LR Statistics For Type 3 Analysis

Source	DF	Chi-Square	Pr > ChiSq
sex	1	10.23	0.0014
ld	1	112.04	<.0001

Compare prob of failure between M & F

$$\log\left(\frac{p}{1-p}\right)_M = \beta_0 + \beta_2 \text{ld} \quad \beta_0 = -2.37$$

$$\log\left(\frac{p}{1-p}\right)_F = \beta_0 + \beta_1 + \beta_2 \text{ld} \quad \beta_2 = 1.54$$

$$\Rightarrow \log \pi_M - \log \pi_F = \log\left(\frac{\pi_M}{\pi_F}\right) = \beta_1$$

where $\pi_M = \frac{p_M}{1-p_M}$, $\pi_F = \frac{p_F}{1-p_F}$ }

$$\Rightarrow \frac{\pi_M}{\pi_F} = e^{\beta_1} \approx 2.72$$

$$\Rightarrow \pi_M = 2.7 \pi_F$$

Odds of failure for M are 2.7 times the odds of failure for F.

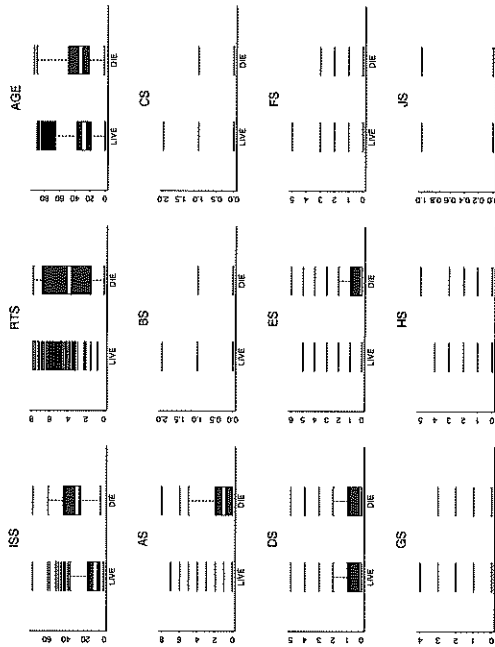
= HEAD, BS = FACE, CS = NECK, DS = THORAX, ES = ABDOMEN, FS = SPINE, GS = UPPER EXTREMITIES, HS = LOWER EXTREMITIES, and JS = SKIN.

I made side-by-side boxplots of the distributions of ISS, AGE, RTS, and AS through JS for the survivors and non-survivors. In several body regions the number of injuries is limited, so these boxplots are not very enlightening. Survivors tend to have lower ISS scores, tend to be slightly younger, tend to have higher RTS scores, and tend to have fewer severe head (AS) and abdomen injuries (ES) than non-survivors. The importance of the effects individually towards predicting survival is directly related to the separation between the survivors and non-survivors scores.

Selecting Predictors in the Trauma Data

The logistic procedure has automated methods for model building, including **backward elimination**, forward selection, and stepwise methods, among others. A SAS proc is given below for performing a backward elimination, starting with a full model having 13 effects: ISS, BP, RTS, AGE, and AS-JS. The elimination process begins by identifying the least significant effect in the full model, based on Wald test p-values. If the least significant effect is not important, the effect is omitted from the model and the resulting model is fitted. The process is repeated until the least significant effect in the model is too important to exclude. This gives the final model. The default significance level for testing effects in backward elimination is .05. The test level was changed to .10 using the `slstay=.10` option on the `model` statement.

In our previous logistic regression analyses, the cases in the data set were pre-aggregated into groups of observations having identical levels of the predictor variables. The numbers of cases in the success category and the group sample sizes were specified in the model statement, along with the names of the predictors. The trauma data set, which is not reproduced here, is **raw data** consisting of one record per patient (i.e. 3132 lines). The logistic model is fitted to data on individual cases by specifying the binary response variable



(SURV) and the predictors on the model statement. Some care is needed because SAS will define the logistic model in such a way that the success category corresponds to the lower of the two levels of response variable. Here SURV=0 corresponds to patients that died whereas SURV=1 corresponds to patients that lived, so the default model is for the log-odds of dying as a function of the predictor variables. To reverse the role of the response levels and model the log-odds of surviving, I used the **descending** option on the model statement (see the printed note from the log window contained with the SAS output). This switch has no impact on the model building process. In particular, you can convert a model for the log-odds of surviving to a model for the log-odds of dying by simply changing the signs of each regression coefficient in the model.

```
data d1;
infile 'C:\Documents and Settings\ed\My Documents\LONGITUDINAL\course\albuq.bed';
input id surv a1 a2 a3 a4 a5 a6 b1 b2 b3 b4 b5 b6 c1 c2 c3 c4 c5 c6
      d1 d2 d3 d4 d5 d6 e1 e2 e3 e4 e5 e6 f1 f2 f3 f4 f5 f6
      g1 g2 g3 g4 g5 g6 h1 h2 h3 h4 h5 h6 j1 j2 j3 j4 j5 j6
      iss iciss bp rts age prob;
as = a3+a4+a5+a6;
bs = b3+b4+b5+b6;
cs = c3+c4+c5+c6;
ds = d3+d4+d5+d6;
es = e3+e4+e5+e6;
fs = f3+f4+f5+f6;
gs = g3+g4+g5+g6;
hs = h3+h4+h5+h6;
js = j3+j4+j5+j6;

proc logistic descending;
model surv = as bs cs ds es fs gs hs js
      iss rts age bp/
      selection=backward scale=none lackfit slistay=.10
      ctable pprob=(.10 to 1 by .10);
run;
```

FROM LOG WINDOW: NOTE: PROC LOGISTIC is modeling the probability that surv=1.

Summary of Backward Elimination					
Step	Effect Removed	DF	Number In	Wald Chi-Square	Pr > ChiSq
1	as	1	12	0.0489	0.8250
2	bs	1	11	0.1196	0.7294
3	fs	1	10	0.1468	0.7016
4	cs	1	9	0.6680	0.4137
5	ds	1	8	0.8078	0.3688

6	gs	1	7	1.1229	0.2893
7	hs	1	6	1.1977	0.2738
8	js	1	5	1.7047	0.1917

Deviance and Pearson Goodness-of-Fit Statistics

Criterion	DF	Value	Value/DF	Pr > ChiSq
Deviance	2573	866.0687	0.3366	1.0000
Pearson	2573	2534.2677	0.9849	0.7031

Number of unique profiles: 2579

Analysis of Maximum Likelihood Estimates

Parameter	DF	Estimate	Standard Error	Chi-Square	Pr > ChiSq
Intercept	1	0.3558	0.4429	0.6454	0.4218
es	1	-0.4613	0.1101	17.5567	<.0001
iss	1	-0.0569	0.00741	58.9887	<.0001
rts	1	0.8431	0.0553	232.1313	<.0001
age	1	-0.0497	0.00529	88.2479	<.0001
bp	1	-0.6351	0.2496	6.4751	0.0109

Odds Ratio Estimates

Effect	Point Estimate	95% Wald Confidence Limits
es	0.630	0.508 0.782
iss	0.945	0.931 0.958
rts	2.324	2.085 2.590
age	0.952	0.942 0.961
bp	0.530	0.325 0.864

I only included the summary table from the backward elimination, and information on the fit of the selected model. The final model includes effects for ES (number of severe abdominal injuries), ISS, RTS, AGE, and BP. All of the effects in the selected model are significant at the 5% level.

Letting p be the probability of survival, the estimated survival probability is given by

$$\log\left(\frac{\hat{p}}{1-\hat{p}}\right) = .3558 - .4613ES - .6351BP - .0569ISS + .8431RTS - .0497AGE.$$

Summary of Logistic Regression Classif.:

Data

	S	D	
S	2825	129	2954
D	40	138	178
	2865	267	3132

Model

$$\text{Correctly Classified} = \frac{2825 + 138}{3132} = 0.946$$

$$\text{False Positive Rate} = \frac{\text{Model} = S \text{ but Data} = D}{\text{Model} = S} = \frac{129}{2954} = 0.044$$

$$\text{False Negative Rate} = \frac{\text{Model} = D \text{ but Data} = S}{\text{Model} = D} = \frac{40}{178} = 0.225$$

$$\text{Sensitivity} = \frac{\text{proportion of survivors correctly classified}}{\text{prop. of survivors}} = \frac{2825}{2825 + 40} = 0.986$$

$$\text{Specificity} = \frac{\text{prop. of non-survivors correctly classified}}{\text{prop. of non-survivors}} = \frac{138}{138 + 129} = 0.517$$

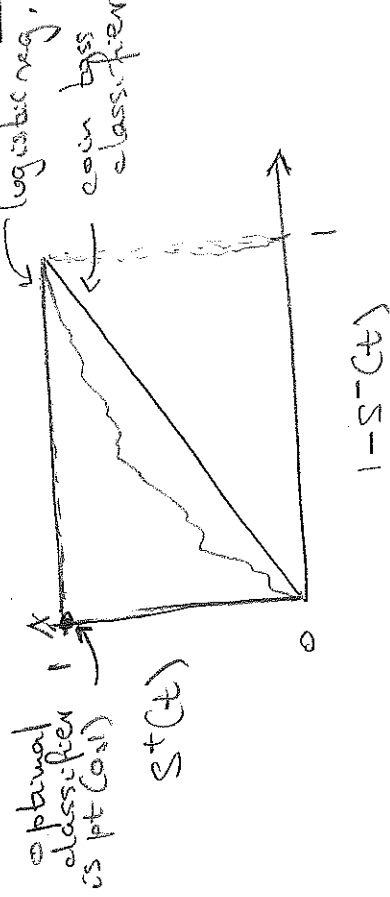
Goals of Classification:

• Simultaneously maximize sensitivity & specificity.

• Sensitivity = $S^+(t)$
Specificity = $S^-(t)$

t = threshold cutoff for classifying (e.g. $p > 0.5$)

• Plot $1 - S^-(t)$ vs. $S^+(t)$ ROC curve



• Summarize ROC curve with AUC = area under curve

$$0.5 \leq \text{AUC} \leq 1.$$

↑
coin toss classifier
↑
optimal classifier