

## The Flight of the Space Shuttle Challenger

On January 28, 1986, the space shuttle Challenger took off on the 25<sup>th</sup> flight in NASA's space shuttle program. Less than 2 minutes into the flight, the spacecraft exploded, killing all on board. A Presidential Commission was formed to explore the reasons for this disaster.

First, a little background information: the space shuttle uses two booster rockets to help lift it into orbit. Each booster rocket consists of several pieces whose joints are sealed with rubber O-rings (0.28 inches wide and 37.5 feet in diameter), which are designed to prevent the release of hot gases produced during combustion. Each booster contains 3 primary O-rings (for a total of 6 for the orbiter). In the 23 previous flights for which there were data (the hardware for one flight was lost at sea), the O-rings were examined for damage.

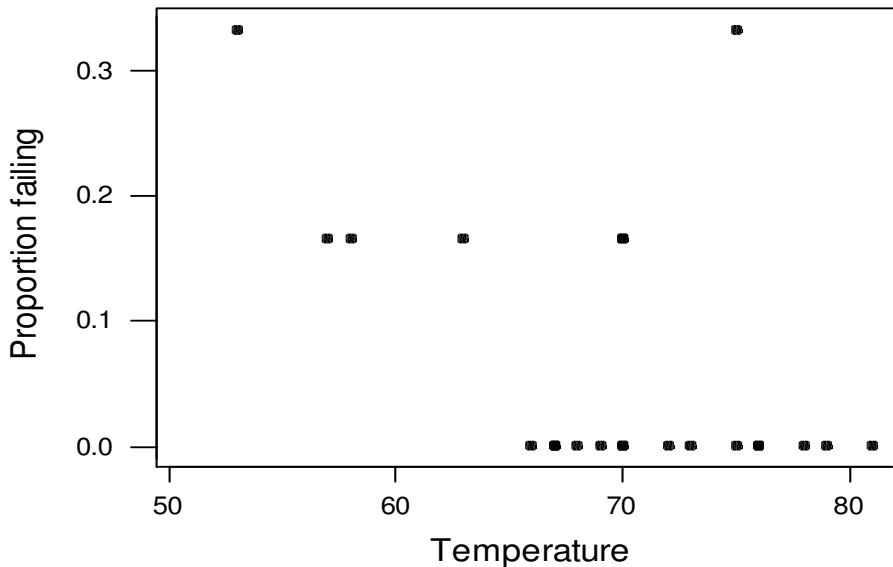
One interesting question is the relationship of O-ring damage to temperature (particularly since it was (forecasted to be) cold — 31° F — on the morning of January 28, 1986). There was a good deal of discussion among the Morton Thiokol engineers the previous day as to whether the flight should go on as planned or not (an important point is that no statisticians were involved in the discussions). A simplified version of one of the arguments made is as follows. There were 7 previous flights where there was damage to at least one O-ring. Consider the following table. The entry  $\hat{p}$  is the frequency estimate of the probability of an O-ring failing for that flight.

Ambient temperature	$\hat{p}$
53°	.333
57°	.167
58°	.167
63°	.167
70°	.167
70°	.167
75°	.333

If you look at the table above, there's no apparent relationship between temperature and the probability of damage; higher damage occurred at both lower and higher temperatures.

Thus, the fact that it was going to be cold on the day of the flight doesn't imply that the flight should be scrubbed. (In fact, this table was not actually constructed the night of January 27th, but was rather given later by two Thiokol staff members as an example of the reasoning in the pre-launch debate. The actual charts faxed from the Thiokol engineers to NASA that night were considerably less informative than even this seriously flawed table.)

Unfortunately, this analysis is completely inappropriate. The problem is that it is ignoring the 16 flights where there was no O-ring damage, acting as if there is no information in those flights. This is clearly absurd! If flights with high temperatures **never** had O-ring damage, for example, that would certainly tell us a lot about the relationship between temperature and O-ring damage! In fact, here is a scatter plot of the frequency estimates of the probability of O-ring damage versus temperature for **all** of the flights:



The picture is very different now. With the exception of the one observation in the upper right of the plot, there is a clear inverse relationship between the probability of O-ring damage and the ambient temperature — lower temperature is associated with higher probability of failure (the unusual observation is the flight of the Challenger from October

30 through November 6, 1985; one way that it was different was that the two O-rings damaged in that flight suffered only “blow-by” [where hot gases rush past the O-ring], while in all of the other flights damaged O-rings suffered “erosion” [where the O-rings burn up], as well as (possibly) blow-by). A plot of this kind would certainly have raised some alarms as to the advisability of launching the shuttle. Unfortunately, such a plot was never constructed.

Here is the full set of data:

Row	Date	Temp	Damaged	O-rings
1	4/12/1981	66	0	6
2	11/12/1981	70	1	6
3	3/22/1982	69	0	6
4	11/11/1982	68	0	6
5	4/4/1983	67	0	6
6	6/18/1983	72	0	6
7	8/30/1983	73	0	6
8	11/28/1983	70	0	6
9	2/3/1984	57	1	6
10	4/6/1984	63	1	6
11	8/30/1984	70	1	6
12	10/5/1984	78	0	6
13	11/8/1984	67	0	6
14	1/24/1985	53	2	6
15	4/12/1985	67	0	6
16	4/29/1985	75	0	6
17	6/17/1985	70	0	6
18	7/29/1985	81	0	6
19	8/27/1985	76	0	6
20	10/3/1985	79	0	6
21	10/30/1985	75	2	6
22	11/26/1985	76	0	6
23	1/12/1986	58	1	6

Logistic regression can be used to analyze the relationship between temperature and the probability of O-ring failure more precisely. In this case, the number of failures is the target variable (which *Minitab* calls *Event*, remember), and the program is told that the number of trials is given in a variable *O-rings* (which is 6 for each flight here). Here is the output of the logistic analysis:

Binary Logistic Regression

Link Function: Logit

Response Information

Variable	Value	Count
Damaged	Event	9
	Non-event	129
O-rings	Total	138

Logistic Regression Table

Predictor	Coef	SE Coef	Z	P	Odds Ratio	95% CI Lower	95% CI Upper
Constant	5.08498	3.05249	1.67	0.096			
Temp	-0.115601	0.0470238	-2.46	0.014	0.89	0.81	0.98

Log-Likelihood = -30.198

Test that all slopes are zero: G = 6.144, DF = 1, P-Value = 0.013

Goodness-of-Fit Tests

Method	Chi-Square	DF	P
Pearson	13.5722	14	0.482
Deviance	11.9564	14	0.610
Hosmer-Lemeshow	5.6769	4	0.225

Table of Observed and Expected Frequencies:

(See Hosmer-Lemeshow Test for the Pearson Chi-Square Statistic)

Value	Group						Total
	1	2	3	4	5	6	
Event							
Obs	0	2	2	0	2	3	9
Exp	0.3	0.6	1.6	1.8	2.0	2.7	
Non-event							
Obs	18	22	34	30	16	9	129
Exp	17.7	23.4	34.4	28.2	16.0	9.3	

Total            18    24    36    30    18    12    138

Measures of Association:

(Between the Response Variable and Predicted Probabilities)

Pairs	Number	Percent	Summary Measures	
Concordant	759	65.4	Somers' D	0.38
Discordant	315	27.1	Goodman-Kruskal Gamma	0.41
Ties	87	7.5	Kendall's Tau-a	0.05
Total	1161	100.0		

The slope coefficient has the following natural interpretation: each increase in temperature by one degree Fahrenheit is associated with an estimated multiplication of the relative odds of an O-ring failure

$$\frac{P(\text{O-ring fails})}{P(\text{O-ring does not fail})}$$

by  $\exp(-.1156) = 0.891$ , or roughly an 11% decrease. This value is given in the output under **Odds Ratio**, along with a 95% confidence interval. If this interval does not contain 1, there is significant predictive power of the predictor on the probability of success (at a .05 level).

There are two other tests given related to the strength of the predictive power of temperature for probability of an O-ring failure. The  $z$ -statistic of  $-2.46$  for **Temperature** corresponds to a  $t$ -statistic in linear regression, and is called a *Wald* statistic (it is equivalent to the odds ratio confidence interval comparison mentioned above). The  $G$ -statistic, given as testing that all slopes are zero, corresponds to the  $F$ -statistic for overall significance in linear regression. Note that for least squares linear regression these two tests are equivalent when there is one predictor, but here the tail probabilities are slightly different, demonstrating that the two tests are not exactly equivalent.

The three goodness-of-fit tests are designed to test whether the logistic model fits the data adequately. All three are based on a  $\chi^2$ -test construction. For each value of temperature given in the data (there are  $J = 16$  distinct values in these data; this is termed the number of distinct *covariate patterns*), let  $\hat{p}_j$  be the fitted probability of O-ring failure,

let  $f_j$  be the observed number of O-rings that failed, and let  $n_j$  be the number of O-rings at risk for that temperature (6 for each flight at that temperature). Note that looking at the data this way means that all flights at a given ambient temperature are pooled together and treated as indistinguishable. These values can be obtained as **Storage** from a logistic fit, and are as follows:

Row	Date	Temp	NOCC1	NTRI1	EPR01
1	4/12/1981	66	0	6	0.072783
2	11/12/1981	70	2	24	0.047106
3	3/22/1982	69	0	6	0.052575
4	11/11/1982	68	0	6	0.058640
5	4/4/1983	67	0	18	0.065357
6	6/18/1983	72	0	6	0.037749
7	8/30/1983	73	0	6	0.033767
8	11/28/1983	70	*	*	*
9	2/3/1984	57	1	6	0.181787
10	4/6/1984	63	1	6	0.099940
11	8/30/1984	70	*	*	*
12	10/5/1984	78	0	6	0.019229
13	11/8/1984	67	*	*	*
14	1/24/1985	53	2	6	0.260787
15	4/12/1985	67	*	*	*
16	4/29/1985	75	2	12	0.026985
17	6/17/1985	70	*	*	*
18	7/29/1985	81	0	6	0.013671
19	8/27/1985	76	0	12	0.024110
20	10/3/1985	79	0	6	0.017166
21	10/30/1985	75	*	*	*
22	11/26/1985	76	*	*	*
23	1/12/1986	58	1	6	0.165220

Note that the  $n_j$  values range from 6 to 24. The sum of the  $n_j$  values is the total sample size, or here 138.

The Pearson goodness-of-fit statistic equals

$$X^2 = \sum_j \frac{(f_j - n_j \hat{p}_j)^2}{n_j \hat{p}_j (1 - \hat{p}_j)},$$

while the deviance statistic equals

$$G^2 = 2 \sum_j \left[ f_j \ln \left( \frac{f_j}{n_j \hat{p}_j} \right) + (n_j - f_j) \ln \left( \frac{n_j - f_j}{n_j (1 - \hat{p}_j)} \right) \right].$$

When the  $n_j$  values are reasonably large, each of these statistics follows a  $\chi^2$  distribution on  $J - p - 1$  degrees of freedom, where  $p$  is the number of predictors in the model, under the null hypothesis that the logistic regression model fits the data. Thus, a small tail probability suggests that the linear logistic regression model is not appropriate for the data. Here both tests have high tail probabilities, indicating no problem with the linear logistic model.

Unfortunately, these tests are not trustworthy when the  $n_j$  values are small (the  $n_j = 6$  values here are marginal). This is the justification for the third goodness-of-fit test, the Hosmer–Lemeshow test. In this test, all of the 138 observations are ordered by estimated O-ring failure probability (of course for these data all of the O-rings for a given flight have the same value of `Temp`, and therefore the same estimated probability of O-ring failure). The observations are then divided into  $g$  roughly equisized groups;  $g$  is usually taken to be 10, except when that would lead to too few observations in each group (as is the case here, where  $g = 6$ ). Based on this new categorization of the data there are values of  $f_j$ ,  $n_j$  and  $n_j \hat{p}_j$ , all of which are given in the Hosmer–Lemeshow table in the output. Then, the Hosmer–Lemeshow goodness-of-fit test is the usual Pearson goodness-of-fit test based on the new categorization, which is compared to a  $\chi^2$  distribution on  $g - 2$  degrees of freedom. It can be seen that the Hosmer–Lemeshow test also does not indicate a lack of fit here. Even the Hosmer–Lemeshow test is suspect, however, when its expected counts for either group are too small (less than two or three, say), which is the case here.

The statistical significance and goodness-of-fit of this model are comforting, of course, but does temperature provide predictive power of any practical importance? Some guidance to answer this question is given in the output under `Measures of Association`. Consider the fitted logistic regression model, with resultant fitted probabilities of O-ring failure  $\hat{p}$  for each of the  $n = 138$  observations. There are  $n_1 = 9$  observed O-ring failures, and  $n_0 = 129$  observed non-failures. Consider each of the pairs  $(i, j)$  of observations where one observation is a failure ( $i$ ) and the other is a non-failure ( $j$ ). There are  $9 \times 129 = 1161$

such pairs, each of which has a corresponding pair  $(\hat{p}_i, \hat{p}_j)$ . We would like the estimated probability of failure to be higher for the observed failure observation than for the observed non-failure observation; that is,  $\hat{p}_i > \hat{p}_j$ . Such pairs are called *concordant*. If for a given pair  $\hat{p}_i < \hat{p}_j$ , the pair is called *discordant*. We would like to have a high percentage of concordant pairs, and a low percentage of discordant pairs. Here there are 65.4% concordant pairs and 27.1% discordant ones, a reasonably good performance. There are no formal cutoffs for what constitutes a “good enough” performance here, but observed values can be compared for different possible models to assess relative practical performance.

The statistics Somers’  $D$ , Goodman–Kruskal  $\gamma$  and Kendall’s  $\tau_a$  are different ways of summarizing these concordancies and discordancies, with higher values indicating more concordancy (e.g.,  $D$  is the difference between concordant and discordant pairs). In particular, Somers’  $D$  is equivalent to the area under the Receiver Operating Characteristic (ROC) curve, and also the Wilcoxon-Mann-Whitney test statistic for comparing the distribution of probability estimates of observations that are successes and failures. Each of these is a measure of the quality of the probability rankings implied by the model (in the sense of concordance), although a good probability ranking need not necessarily be *well-calibrated*. That is, if the estimated probability of success for each observation was exactly one-half the true probability, the probability rankings would be perfect, but not well-calibrated, since the estimates are far from the true probabilities.

In the medical diagnostic testing literature, the following rough guide for interpretation of  $D$  has been suggested; it is perhaps useful as a way to get a sense of what the value is telling you, but should be recognized as being fairly arbitrary, and should not be taken overly seriously.

Range of $D$	Rough interpretation
0.8 – 1.0	Excellent separation
0.6 – 0.8	Good separation
0.4 – 0.6	Fair separation
0.2 – 0.4	Poor separation
0.0 – 0.2	Little to no separation

Just as is true for other regression models, unusual observations can have a strong effect on a fitted logistic regression model. Among the diagnostics that are available for logistic regression are three that roughly correspond to the standardized residuals (here the standardized Pearson residuals), Cook's distance (here the standardized Delta-beta  $[\Delta\beta]$ ) and leverage values. In this context the appropriate guidelines for the leverage values is  $(2.5)(p + 1)/J$ , where  $J$  is the number of distinct covariate patterns. Here are the values here:

Row	SPRE1	DSBE1	HI1
1	-0.70527	0.02790	0.053119
2	0.94160	0.23361	0.208538
3	-0.59239	0.01894	0.051209
4	-0.62752	0.02110	0.050853
5	-1.21983	0.27107	0.154100
6	-0.49902	0.01443	0.054772
7	-0.47134	0.01322	0.056146
8	*	*	*
9	-0.10892	0.00340	0.222683
10	0.56529	0.02428	0.070609
11	*	*	*
12	-0.35377	0.00799	0.060045
13	*	*	*
14	0.56258	0.29503	0.482442
15	*	*	*
16	3.17765	1.33690	0.116919
17	*	*	*
18	-0.29729	0.00555	0.059061
19	-0.57995	0.04523	0.118540
20	-0.33389	0.00712	0.059995
21	*	*	*
22	*	*	*
23	0.01054	0.00002	0.180967

There is an apparent outlier at row 16, corresponding to an ambient temperature of  $75^\circ$ . Unfortunately, since there are two flights at that temperature, we can't tell for sure which is actually the outlier (of course, in this case we know what it is from the earlier

graph, but in general the collapsing approach of Minitab makes it difficult to tell which observation is actually an outlier if there are replications in the data).

For this reason, it is appropriate to remove the collapsing effect by forcing each observation to have a unique set of predictor variable values, at least when looking at diagnostics. The way this is done is by “jittering” at least one predicting variable by adding a small amount of random noise to the variable. In fact, in this context we know that the temperature values are only given to the nearest integer value, so this is not at all unreasonable, but even if the values were exact, we need to do this if we want to examine diagnostics. Note that this is only sensible for truly numerical variables; it is not appropriate to jitter a categorical predictor, including a 0/1 variable. Here are the results of a logistic fit using jittered temperature:

Binary Logistic Regression: Damaged, O-rings versus Temperature

Link Function: Logit

Response Information

Variable	Value	Count
Damaged	Event	9
	Non-event	129
O-rings	Total	138

Logistic Regression Table

Predictor	Coef	SE Coef	Z	P	Odds	95% CI	
					Ratio	Lower	Upper
Constant	5.08495	3.05249	1.67	0.096			
Temperature	-0.115601	0.0470238	-2.46	0.014	0.89	0.81	0.98

Log-Likelihood = -30.198

Test that all slopes are zero: G = 6.144, DF = 1, P-Value = 0.013

Goodness-of-Fit Tests

Method	Chi-Square	DF	P
--------	------------	----	---

Pearson	29.9804	21	0.092
Deviance	18.0864	21	0.644
Hosmer-Lemeshow	7.7170	6	0.260

Table of Observed and Expected Frequencies:  
(See Hosmer-Lemeshow Test for the Pearson Chi-Square Statistic)

Value	Group								Total
	1	2	3	4	5	6	7	8	
Event									
Obs	0	2	0	1	1	0	2	3	9
Exp	0.3	0.5	0.6	0.8	0.9	1.2	2.0	2.7	
Non-event									
Obs	18	16	18	17	17	18	16	9	129
Exp	17.7	17.5	17.4	17.2	17.1	16.8	16.0	9.3	
Total	18	18	18	18	18	18	18	12	138

Measures of Association:  
(Between the Response Variable and Predicted Probabilities)

Pairs	Number	Percent	Summary Measures
Concordant	759	65.4	Somers' D 0.38
Discordant	315	27.1	Goodman-Kruskal Gamma 0.41
Ties	87	7.5	Kendall's Tau-a 0.05
Total	1161	100.0	

As would be expected, the output changes very little, with the exception of the Pearson and Deviance statistics (recall that their construction depends on how covariate patterns are defined). Here are the diagnostics; note that each flight now has a distinct set of diagnostics, which is why you **must** jitter the data if collapsing occurs because of non-unique covariate patterns based on numerical predictor(s):

Row	SPRE2	DSBE2	HI2	EPR02
1	-0.70527	0.02790	0.053119	0.072784
2	1.41982	0.11088	0.052135	0.047106
3	-0.59239	0.01894	0.051209	0.052575
4	-0.62752	0.02110	0.050853	0.058641
5	-0.66505	0.02395	0.051367	0.065358
6	-0.49902	0.01443	0.054771	0.037750

7	-0.47134	0.01322	0.056146	0.033767
8	-0.55939	0.01721	0.052135	0.047106
9	-0.10892	0.00340	0.222683	0.181786
10	0.56530	0.02428	0.070609	0.099939
11	1.41980	0.11087	0.052134	0.047106
12	-0.35377	0.00799	0.060045	0.019229
13	-0.66504	0.02395	0.051367	0.065356
14	0.56259	0.29503	0.482443	0.260786
15	-0.66505	0.02395	0.051367	0.065358
16	-0.42040	0.01097	0.058459	0.026986
17	-0.55939	0.01721	0.052135	0.047106
18	-0.29729	0.00555	0.059061	0.013671
19	-0.39696	0.00993	0.059270	0.024110
20	-0.33389	0.00712	0.059995	0.017166
21	4.77254	1.41422	0.058460	0.026985
22	-0.39696	0.00993	0.059270	0.024111
23	0.01055	0.00002	0.180967	0.165219

Now it's clear that the previously mentioned flight (number 21) is a very clear outlier, with 2 of 6 O-rings damaged when the estimated probability of O-ring damage was only .027. Here is output from the data with that flight omitted (I'm sticking with the jittered data):

#### Binary Logistic Regression

Link Function: Logit

#### Response Information

Variable	Value	Count
Damaged	Event	7
	Non-event	125
O-rings	Total	132

#### Logistic Regression Table

Predictor	Coef	SE Coef	Z	P	Odds	95% CI	
					Ratio	Lower	Upper
Constant	8.66156	3.63428	2.38	0.017			
Temperature	-0.176805	0.0586846	-3.01	0.003	0.84	0.75	0.94

Log-Likelihood = -22.041

Test that all slopes are zero:  $G = 10.657$ ,  $DF = 1$ ,  $P\text{-Value} = 0.001$

#### Goodness-of-Fit Tests

Method	Chi-Square	DF	P
Pearson	13.4054	20	0.859
Deviance	9.4097	20	0.978
Hosmer-Lemeshow	7.9734	6	0.240

#### Table of Observed and Expected Frequencies:

(See Hosmer-Lemeshow Test for the Pearson Chi-Square Statistic)

Value Event	Group								Total
	1	2	3	4	5	6	7	8	
Obs	0	0	0	2	0	0	3	2	7
Exp	0.1	0.2	0.3	0.4	0.6	0.8	2.6	2.0	
Non-event									
Obs	18	18	18	16	18	18	15	4	125
Exp	17.9	17.8	17.7	17.6	17.4	17.2	15.4	4.0	
Total	18	18	18	18	18	18	18	6	132

#### Measures of Association:

(Between the Response Variable and Predicted Probabilities)

Pairs	Number	Percent	Summary Measures
Concordant	671	76.7	Somers' D 0.61
Discordant	137	15.7	Goodman-Kruskal Gamma 0.66
Ties	67	7.7	Kendall's Tau-a 0.06
Total	875	100.0	

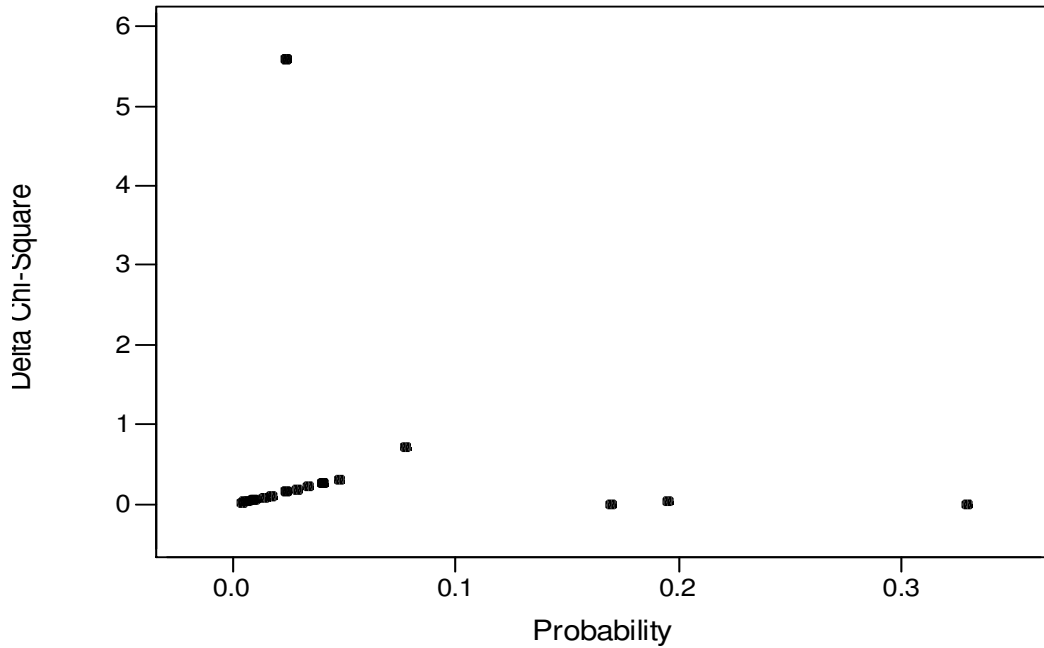
The strength of the relationship has gone up considerably once the outlier is removed, with there now being an estimated 16% reduction in the odds of an O-ring being damaged with each additional degree of temperature at launch. The goodness-of-fit tests suggest

no lack of fit (remember, the Pearson and deviance tests are at least marginally valid here, since there are 6 replications for each flight). Here are diagnostics:

Row	Temp	SPRE3	DSBE3	HI3	EPR03
1	66	-0.56300	0.021910	0.064653	0.047087
2	70	2.36708	0.346866	0.058297	0.023781
3	69	-0.43075	0.011820	0.059890	0.028251
4	68	-0.47093	0.014497	0.061358	0.033532
5	67	-0.51488	0.017782	0.062859	0.039759
6	72	-0.32946	0.006253	0.054467	0.016818
7	73	-0.30122	0.004998	0.052208	0.014131
8	70	-0.39397	0.009609	0.058297	0.023781
9	57	-0.20096	0.011923	0.227940	0.195241
10	63	0.85043	0.060028	0.076638	0.077475
11	70	2.36705	0.346858	0.058298	0.023782
12	78	-0.19226	0.001493	0.038827	0.005886
13	67	-0.51487	0.017781	0.062859	0.039758
14	53	0.02772	0.000975	0.559114	0.329800
15	67	-0.51488	0.017782	0.062859	0.039759
16	75	-0.25174	0.003135	0.047138	0.009964
17	70	-0.39397	0.009609	0.058297	0.023781
18	81	-0.14685	0.000682	0.030665	0.003472
19	76	-0.23011	0.002461	0.044416	0.008363
20	79	-0.17574	0.001155	0.036043	0.004937
21	76	-0.23011	0.002461	0.044416	0.008363
22	58	-0.01646	0.000060	0.180323	0.168946

There are no extreme outliers, but the low temperature cases are possible leverage points (this is not surprising, given that most launches were at temperatures over  $65^\circ$ ). The noteworthy  $70^\circ$  observations correspond to two  $70^\circ$  flights where there was an O-ring failure. Omitting these two flights doesn't change things very much (strengthening the relationship further), although a plot of the change in the Pearson statistic versus estimated probability does show the two points as unusual.

### Delta Chi-Square versus Probability



Binary Logistic Regression

Link Function: Logit

Response Information

Variable	Value	Count
Damaged	Event	5
	Non-event	115
O-rings	Total	120

Logistic Regression Table

Predictor	Coef	SE Coef	Z	P	Odds	95% CI	
					Ratio	Lower	Upper
Constant	14.4116	5.70596	2.53	0.012			
Temperature	-0.280180	0.0985633	-2.84	0.004	0.76	0.62	0.92

Log-Likelihood = -13.320

Test that all slopes are zero: G = 14.929, DF = 1, P-Value = 0.000

Goodness-of-Fit Tests

Method	Chi-Square	DF	P
Pearson	3.44529	18	1.000
Deviance	2.78174	18	1.000
Hosmer-Lemeshow	1.56885	8	0.992

Table of Observed and Expected Frequencies:

(See Hosmer-Lemeshow Test for the Pearson Chi-Square Statistic)

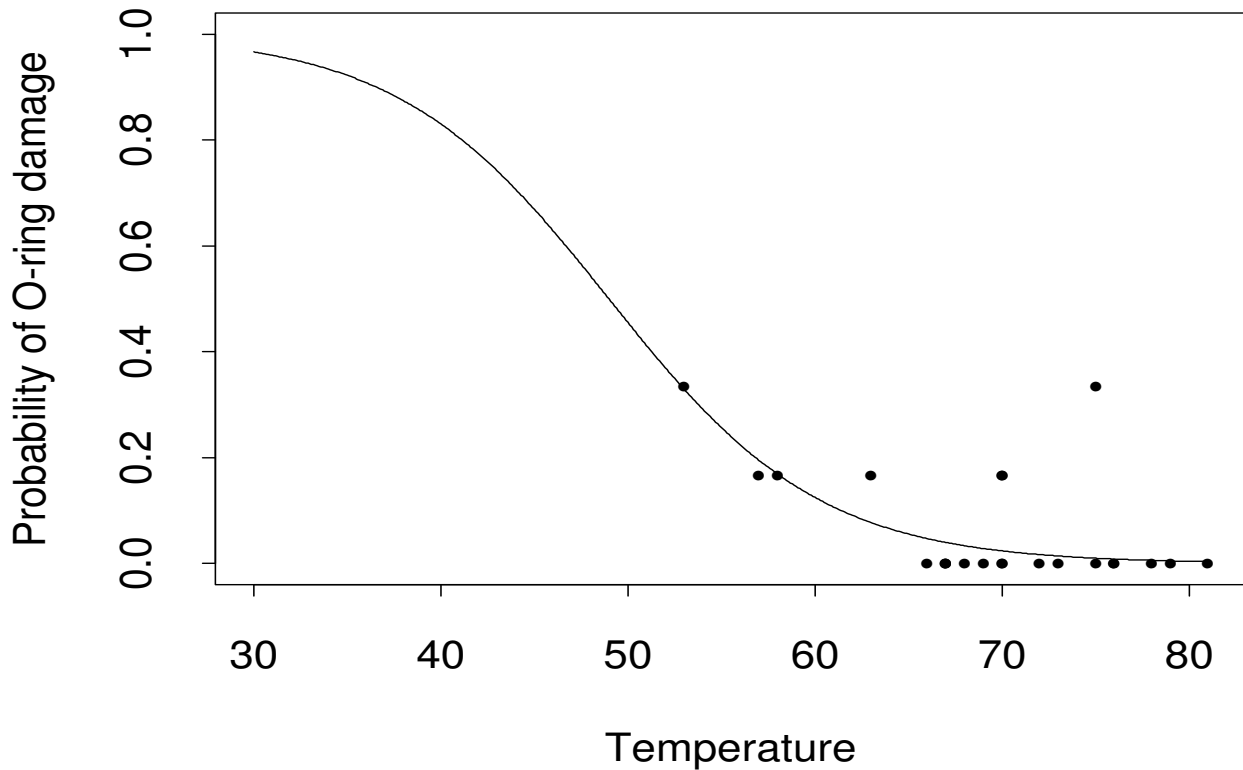
Value	Group										Total	
	1	2	3	4	5	6	7	8	9	10		
Event												
Obs	0	0	0	0	0	0	0	0	2	3		5
Exp	0.0	0.0	0.0	0.0	0.1	0.1	0.2	0.2	1.0	3.4		
Non-event												
Obs	12	12	12	12	12	12	12	12	10	9		115
Exp	12.0	12.0	12.0	12.0	11.9	11.9	11.8	11.8	11.0	8.6		
Total	12	12	12	12	12	12	12	12	12	12	12	120

Measures of Association:

(Between the Response Variable and Predicted Probabilities)

Pairs	Number	Percent	Summary Measures
Concordant	525	91.3	Somers' D 0.87
Discordant	27	4.7	Goodman-Kruskal Gamma 0.90
Ties	23	4.0	Kendall's Tau-a 0.07
Total	575	100.0	

What about the morning of January 28, 1986? Here is a plot of the logistic curve for different values of temperature based on all flights except the October/November 1985 flight:



Substituting into the logistic function gives a probability estimate of O-ring failure for a temperature of 31° of **.96!** (This is an extrapolation, but you get the idea.) Indeed, with the benefit of hindsight, it can be seen that the Challenger disaster was not at all surprising, **given data that were available at the time of the flight**. As a result of its investigations, one of the recommendations of the commission was that a statistician be part of the ground control team from that time on. A complete (and more correct) discussion of this material can be found in the paper “Risk Analysis of Space Shuttle: Pre-Challenger Prediction of Failure,” by S.R. Dalal, E.B. Fowlkes and B.A. Hoadley, *Journal of the American Statistical Association*, **84**, 945–957 (1989). Chapter 2 of Edward R. Tufte’s 1997 book *Visual Explanations: Images and Quantities, Evidence and Narrative* discusses the background of the disaster, and the charts used by the Thiokol engineers in their discussions with NASA.

By the way, an alternative way that these data might have been presented was as a set of 138 observations (one for each O-ring, rather than one for each flight), with a 0/1 target variable reflecting failure or non-failure of each O-ring. This is what the dataset would look like (now ordered by temperature):

Row	Temp	Failed
1	53	1
2	53	1
3	53	0
4	53	0
5	53	0
6	53	0
7	57	1
8	57	0
9	57	0
10	57	0
11	57	0
12	57	0
13	58	1
14	58	0
15	58	0
16	58	0
17	58	0
18	58	0
19	63	1
20	63	0
21	63	0
22	63	0
23	63	0
24	63	0
25	66	0
26	66	0
27	66	0
28	66	0
29	66	0
30	66	0
31	67	0
32	67	0
33	67	0
34	67	0
35	67	0
36	67	0
37	67	0
38	67	0
39	67	0
40	67	0
41	67	0

42	67	0
43	67	0
44	67	0
45	67	0
46	67	0
47	67	0
48	67	0
49	68	0
50	68	0
51	68	0
52	68	0
53	68	0
54	68	0
55	69	0
56	69	0
57	69	0
58	69	0
59	69	0
60	69	0
61	70	0
62	70	0
63	70	0
64	70	0
65	70	0
66	70	0
67	70	0
68	70	0
69	70	0
70	70	0
71	70	0
72	70	0
73	70	1
74	70	0
75	70	0
76	70	0
77	70	0
78	70	0
79	70	1
80	70	0
81	70	0
82	70	0
83	70	0
84	70	0
85	72	0

86	72	0
87	72	0
88	72	0
89	72	0
90	72	0
91	73	0
92	73	0
93	73	0
94	73	0
95	73	0
96	73	0
97	75	0
98	75	0
99	75	0
100	75	0
101	75	0
102	75	0
103	75	1
104	75	1
105	75	0
106	75	0
107	75	0
108	75	0
109	76	0
110	76	0
111	76	0
112	76	0
113	76	0
114	76	0
115	76	0
116	76	0
117	76	0
118	76	0
119	76	0
120	76	0
121	78	0
122	78	0
123	78	0
124	78	0
125	78	0
126	78	0
127	79	0
128	79	0
129	79	0

130	79	0
131	79	0
132	79	0
133	81	0
134	81	0
135	81	0
136	81	0
137	81	0
138	81	0

Which representation is better? It turns out not to matter; if you analyze the data in this form, where **Failed** is chosen as the **Response** variable in the Minitab dialog box, the resultant output will be identical to that obtained using the data represented at the level of 23 different flights. There is one advantage to the earlier representation, however; since the natural way to view these data is at the flight level, rather than the O-ring level, jittering the data in the flight-level form is more natural.

### Minitab commands

Logistic regression modeling is obtained by clicking on **Stat** → **Regression** → **Binary Logistic Regression**. There are various ways that the data might be presented, which affect the command structure to the program. The two most common forms are as follows:

- (1) The target variable is given as the number of successes out of the number of trials (or the number of items “at risk”). Click the radio button next to **Response in event/trial format**, enter the variable with the number of successes in the box next to **Number of events:**, and enter the variable with the number of trials in the box next to **Number of trials:**.
- (2) The target variable is a 0/1 variable that represents success or failure for each observation. Enter the name of this variable in the box next to **Response:**.

The predicting variables for the model are entered under **Model:**. This includes both continuous variables and categorical ones. Categorical variables must also be entered under **Factors (optional):**. Interactions with (and between) factors are entered under **Model:** using the “multiplication” form as in ANOVA and ANCOVA modeling.

Diagnostics, such as standardized Pearson residuals, leverages, and (standardized) delta betas are obtained by clicking on **Storage**. This dialog box also allows storage of fitted success probabilities under **Event probability**, the number of successes for each distinct covariate pattern under **Number of occurrences of the event**, and the number of trials for each covariate pattern under **Number of trials**. Diagnostic plots, such as one of Delta Chi-Square versus Event Probability, are obtained by clicking on **Graphs**.

To jitter a predicting variable, follow these steps:

- (1) Go to **Calc** → **Random Data** → **Uniform**. In the box labeled **Number of rows of data to generate**, enter the sample size (the number of observations, not the number of covariate patterns). Put in a new variable name under **Store in columns:**, such as **Jitter**. Under **Lower endpoint:** enter a negative number close to zero, such as  $-.001$ . Under **Upper endpoint:** enter the positive version of this number ( $.001$ , for example). The numbers you choose here should be smaller than the resolution of all of your predictors; so, for example, if one predictor is given to three decimal digits, use  $-.0001$  and  $.0001$  here.
- (2) Go to the calculator, and create new versions of each predictor you wish to jitter as sums of the predictor with **Jitter**. So, for example, the variable **Newpred1** is determined as **Predict1 + Jitter**.
- (3) Fit the logistic regression model using the new predictors in place of the old ones.