1. The following table contains data from a study concerning the distribution of party aliation in a neighborhood. The interest was whether there was an association between registered political parties and the race.

Party	White	Black	Hispanic	Asian	Total
Democrat	431	190	180	95	896
Republican	390	248	159	106	902
Independent	406	177	164	82	828
Total	1227	615	503	282	2627

To determine whether there was an association between registered political party and the race, we would like to perform the Chi-square test.

a. We first calculate the table of expected counts. Table below lists the expected counts calculated in SAS. The table is incomplete as results for four cells in the lower-right corner of the table are removed and replaced by a; b; c; d respectively. Complete the table by calculating the values of a; b; c; d. Remember to show your calculation process.

Party	White	Black	Hispanic	Asian	Total
Democrat	a	b	171.71	96.19	896
Republican	c	d	172.96	96.89	902
Independent	386.75	193.81	158.73	88.92	828
Total	1227	615	503	282	2627

d = 633.34 - 632.15 = 1.19 b = 421.19 - 1.19 = 420

$$c + 1.19 = 632.15$$
; $c = 632.15 - 1.19 = 630.96$

In summary:

a = 209.29

b = 420

c = 630.96

d = 1.19

b. What are the null hypothesis and alternative hypothesis of this test?

 H_{0} : There is no statistically significant relationship between registered political parties and race.

 H_{A} : There is a statistically significant relationship between registered political parties and race.

c. Performing data analysis in SAS, we got the chi-square test statistic χ^2 = 16.50 (with p-value 0.011). What are the degrees of freedom for this χ^2 test statistic?

$$df = (r - 1)(c - 1)$$

$$df = (3 - 1)(4 - 1)$$

$$df = (2)(3) = 6$$

There are 6 degrees of freedom for this chi-square test.

2. A survey is conducted to see whether people have a certain symptom X. The following table breaks down these results by gender:

Symptom X by Gender

having symptom X	Male	Female
Yes	9	36
No	39	63

a. Under the hypothesis of independence between row and column variables, calculate the expected cell count for the two cells in the column for women.

ECC = (row i total x column j total) / total count in table

$$ECC = (36) * (36 + 63) / (9 + 36 + 39 + 63)$$

$$ECC = (36 * 99) / 147$$

b. What is the chi-square χ^2 statistic for this contingency table? Is it significant at the α = 0.05 level? Justify your answer. (Use Chi-square distribution table in textbook)

```
1  /* Re-create contingency table */
2  data symptomx;
3   input symptomX gender$ count;
4  datalines;
5  1  Male 9
6  0  Male 39
7  1  Female 36
8  0  Female 63
9 ;
10  run;
11
12  proc freq data = symptomx;
   tables symptomX*gender / chisq expected nocol norow nopercent;
   weight count;
15  run;
```

	Table of symptomX by gender					
Expected			gender			
sym	ptomX	Female	Male	Total		
	0	63 68.694				
	1	36 30.306	_			
Tota	ıl	99	48	147		
Statistics for Ta	able of s	ymptom	C by gend	ler		
Statistics for Ta	able of s	ymptom)	C by gend	der Prob		
	able of s					
Statistic		DF 1	Value	Prob		
Statistic Chi-Square	i-Square	DF 1	Value 4.7215	Prob 0.0298		
Statistic Chi-Square Likelihood Ratio Ch	i-Square Square	DF 1 1 1 1	Value 4.7215 4.9803	Prob 0.0298 0.0256		
Statistic Chi-Square Likelihood Ratio Ch Continuity Adj. Chi-	i-Square Square	DF 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	Value 4.7215 4.9803 3.9287	Prob 0.0298 0.0256 0.0475		
Statistic Chi-Square Likelihood Ratio Ch Continuity Adj. Chi- Mantel-Haenszel Ch	i-Square Square i-Square	DF 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1	Value 4.7215 4.9803 3.9287 4.6894	Prob 0.0298 0.0256 0.0475		

Based on the SAS output, the χ^2 statistic for this table is 4.7215 with a p-value of 0.0298. df = (2 - 1)(2 - 1) = (1)(1) = 1. The critical value for $\alpha = 0.05$ and 1 degree of freedom is 3.841. Because the χ^2 statistic is greater than the critical value, it is statistically significant, and gives us the evidence to reject the null hypothesis.

c. Calculate the relative odds (i.e., the OR) of having symptom X, for men to women.

```
OR^* = ad / bc = (9 * 63) / (36 * 39) = 567 / 1404 = 0.4038
```

 $\exp(-1.7392), \exp(-0.0744) = (0.1757, 0.9283)$

d. Calculate 95% confidence interval for the relative odds (i.e., the OR) in part (c). $\log(\text{OR}^{\,\circ}) = \ln(0.4038) = -0.9068$ $\operatorname{se}(\log(\text{OR}^{\,\circ})) = \sqrt{[(1/9) + (1/63) + (1/36) + (1/39)]} = 0.4247$ Z(95%) = 1.96 $\operatorname{Lower} = -0.9068 - 1.96 * 0.4247 = -1.7392$ $\operatorname{Upper} = -0.9068 + 1.96 * 0.4247 = -0.0744$

- 3. The dataset lowbwt.sas7bdat contains information for the sample of 100 low birth weight infants born in Boston, Massachusetts. The variable grmhem is a dichotomous random variable indicating whether an infant experienced a germinal matrix hemorrhage. The value 1 indicates that a hemorrhage occurred and 0 that it did not. The infants' five-minute apgar scores are saved under the name apgar5, and indicators of toxemia where 1 represents a diagnosis of toxemia during pregnancy for the child's mother and 0 no such diagnosis under the variable name tox.
 - a. Write down the equation for a logistic regression model where germinal matrix hemorrhage is the response and five-minute apgar score is the predictor, using β_1 to represent the regression coefficient of apgar score. log{p / (1 p)} = β_0 + β_1 apgar5 Where p represents the grmhem variable, β_0 represents the intercept, β_1 represents the regression coefficient of the apgar score, and apgar5 represents the predictor variable apgar score.
 - b. Fit the logistic regression model in part (a). What is β_1 , the estimated regression coefficient of apgar score? What's the interpretation of β_1 ?

```
libname ADS534 '/home/u63483466/ADS534/Data';

data work.lowbwt;

set ADS534.lowbwt;

run;

/* Fit the logistic regression model */
proc logistic data = work.lowbwt;

model grmhem = apgar5;
run;
```

	Analy	sis of Maxi	mum Lik	kelihood E	stimates	3
Parameter	DF	Estimate	Standa En		Wald Square	Pr > ChiSq
Intercept	1	0.3037	0.61	191	0.2407	0.6237
apgar5	1	0.2496	0.10)44	5.7206	0.0168
		Odds	Ratio Es	stimates		
	Effect	Point Es	timate	95% Confider	ts	
	apgar:	5	1.284	1.046 1.575		' 5

The value of β_1 is 0.2496. This means that for every 1-unit increase in the five-minute apgar score, the likelihood of grmhem increases by 0.2496. While grmhem has a value of 1 or 0, this means that higher apgar scores increase the likelihood of a germinal matrix hemorrhage changing from 0 to 1, or occurring.

c. If a particular child has a five-minute apgar score of 3, what is the predicted probability that this child will experience a brain hemorrhage?

```
27 /* Predict the probability for Apgar score of 3 */
28 /* Calculate predicted probabilities */
29 proc logistic data = work.lowbwt;
30 model grmhem = apgar5 / LINK = LOGIT;
     output out = PredictedProb1 predicted = p_GrmHem;
31
32 run;
33
34 /* Calculate predicted probability for apgar5 = 3 */
35 data PredictedProb1;
     set PredictedProb1;
37
     if apgar5 = 3;
38 run;
39
40 /* Display predicted probability */
41 proc print data = PredictedProb1 noobs;
    var apgar5 p_GrmHem;
43 run;
 apgar5 p GrmHem
    3
        0.74126
```

With an apgar5 score of 3, the probability of experiencing a germinal matrix hemorrhage is 0.74126, or 74.126%.

d. What is the estimated odds ratio of suffering a germinal matrix hemorrhage associated with 1 unit increase in five-minute apgar score?

```
/* Determine the OR of germ with 1-unit increase */
proc logistic data = work.lowbwt;
model grmhem = apgar5 / LINK = LOGIT;
run;
```

Odds Ratio Estimates					
Effect	Point Estimate	95% Confiden			
apgar5	1.284	1.046	1.575		

Based on the output from SAS, a 1-unit increase in apgar5 score will increase the odds of suffering a germinal matrix hemorrhage by 1.284, and the 95% confidence limits suggest that the value will most likely fall between 1.046 and 1.575.

e. What is the estimated odds ratio of suffering a germinal matrix hemorrhage associated with 3 units increase in five-minute apgar score?

```
27 /* Predict the probability for Apgar score of 3 */
28 /* Calculate predicted probabilities */
29 proc logistic data = work.lowbwt;
      model grmhem = apgar5 / LINK = LOGIT;
31
      output out = PredictedProb1 predicted = p_GrmHem;
32 run;
33
34 | /* Calculate predicted probability for apgar5 = 3 */
35 data PredictedProb1;
       set PredictedProb1;
37
      if apgar5 = 3;
38 run;
39
40 /* Display predicted probability */
41 proc print data = PredictedProb1 noobs;
       var apgar5 p_GrmHem;
43 run;
         Odds Ratio Estimates
                    95% Wald
  Effect | Point Estimate | Confidence Limits
             1 284 1 046
                          1 575
  apgar5
```

Based on the SAS output, the estimated odds ratio associated with a 3-unit increase in apgar5 is 1.284, with 95% confidence limits of (1.046, 1.575).

f. Write down the equation for a logistic regression model where germinal matrix hemorrhage is the response and toxemia status is the predictor, using β_2 to represent the regression coefficient of toxemia status.

$$log{p / (1 - p)} = \beta_0 + \beta_2 tox$$

In this equation, p represents the grmhem variable, β_0 represents the intercept, and β_2 tox represents the change in grmhem given a 1-unit increase in toxemia status.

g. Fit the logistic regression model in part (f) in SAS. What is β_2^* , the estimated regression coefficient of toxemia status? What's the interpretation of β_2^* ?

```
/* Fit the logistic regression model */
proc logistic data = work.lowbwt;
model grmhem = tox / LINK = LOGIT;
run;
```

	Analy	sis of Maxii	mum Li	kelih	ood E	stimates	•
Parameter	DF	Estimate	Stand Er	ard rror	Chi-	Wald Square	Pr > ChiSq
Intercept	1	1.5353	0.2	946	2	7.1530	<.0001
tox	1	1.4604	1.0662		1.8761		0.1708
		Odds	Ratio E	stima	tes		
	Effect	Point Est	timate	Con	95% fiden	Wald ce Limit	S

 β°_{2} , the estimated regression coefficient of toxemia status, is 1.4604. This indicates that an increase from 0 to 1 for toxemia will result in a 1.4604 increase

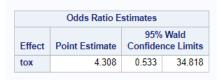
in grmhem. Given that grmhem is also on a scale from 0 to 1, this means that higher tox scores will likely result in higher odds of a germinal matrix hemorrhage occurring.

h. For a child whose mother was diagnosed with toxemia during pregnancy, what is the predicted probability of experiencing a germinal matrix hemorrhage?

```
43 /* Calculate predicted probabilities */
44 proc logistic data = work.lowbwt;
      model grmhem = tox / LINK = LOGIT;
     output out = PredictedProb predicted = p_GrmHem;
47 run;
48
49 /* Calculate predicted probability for toxemia status = 1 */
50 data PredictedProb;
51 set PredictedProb;
    if tox = 1;
52
53 run;
54
55 /* Display predicted probability */
56 proc print data = PredictedProb noobs;
    var tox p_GrmHem;
57
58 run;
  tox p_GrmHem
   1
        0.95238
```

Based on a toxemia status of '1', the predicted probability is 0.95238, or 95.238%. This indicates a high likelihood of a germinal matrix hemorrhage occurring with a diagnosis of toxemia.

i. What are the estimated odds of suffering a germinal matrix hemorrhage for children whose mothers were diagnosed with toxemia relative to children whose mothers were not?



The odds ratio suggests that a change in toxemia status is met with odds being 4.308 times more likely to suffer a germinal matrix hemorrhage compared to no toxemia being present.