

MSPA PREDICT 411

Bonus Problem: Chapter 2

Introduction

This document presents the results of first set of bonus problems for the Masters of Science in Predictive Analytics course: PREDICT 411. This assessment required the student to work through the problem set of Chapters 1 & 2 of Hoffmann (2004), Generalized Linear Models, An Applied Approach.

Question 1&2

Specify the probability distributions that best describe the following variables. Suppose you wish to analyze each of these variables using regression techniques. Select the most likely link function for each distribution.

Part A

A measure of the number of avalanches that occur per year in the Wasatch mountain range of Utah

Poisson, with link function: $X\beta = \ln(\mu)$

Part B

A measure of whether or not members of a large, nationally representative sample of adults smoke cigarettes.

Binomial, with link function: $X\beta = \ln\left(\frac{\mu}{1-\mu}\right)$

Part C

A measure of the temperature (in Kelvin) inside a sample of volcanoes in Japan.

Inverse Gaussian, with link function: $X\beta = -\mu^{-2}$

Part D

A measure of whether members of a sample have done one of the following mutually exclusive events in the past year: Remained with their religious denomination, joined different religious denomination, or left their religious denomination without joining another.

Multinomial, with link function: $X\beta = \ln\left(\frac{\mu}{1-\mu}\right)$

Part E

A measure of whether or not firms in a national registry have adopted a public venture capital program.

Binomial, with link function: $X\beta = \ln\left(\frac{\mu}{1-\mu}\right)$

Part F

A measure of whether members of a sample of workers have either quit a job, been laid-off from a job, been fired from a job, or remained in their jobs in the past year.

Multinomial, with link function: $X\beta = \ln\left(\frac{\mu}{1-\mu}\right)$

Part G

In a sample of adult probationers in Oregon, a measure of the number of times arrested in the previous ten years.

Poisson, with link function: $X\beta = \ln(\mu)$

Question 3

Compute the expected values (means) and variances for each of the following variables.

Part A

A sample of 1,500 adults in which the probability of alcohol use is 0.65.

Using binomial distribution, $E(X) = n \times p$, $Var(X) = n \times p(1 - p)$

expected value: $1500 \times 0.65 = 975$, variance: $1500 \times 0.65(1 - 0.65) = 341.25$

Part B

A sample of 200 adults with the following probabilities of involvement in the workforce: 0.55 of being employed full-time, 0.15 of being employed part-time, 0.10 of being unemployed, and 0.20 of not participating in the workforce (e.g. homemakers, students).

Using multinomial distribution, $E(X) = p_1 \times n$ for each group (denoted as 1 for first group). $Var(X) = n \times p_1(1 - p_1)$ for each group (denoted as 1 for first group).

- 0.55 full time, expected value: $0.55 \times 200 = 110$, variance: $200 \times 0.55(1 - 0.55) = 49.5$
- 0.15 part time, expected value: $0.15 \times 200 = 30$, variance: $200 \times 0.15(1 - 0.15) = 25.5$
- 0.10 unemployed, expected value: $0.10 \times 200 = 20$, variance: $200 \times 0.10(1 - 0.10) = 18$
- 0.20 not participating, expected value: $0.20 \times 200 = 40$, variance: $200 \times 0.20(1 - 0.20) = 32$

Part C

A sample of 850 adolescents with the following probabilities of low and high self-esteem: 0.45 low self-esteem; and 0.55, high self-esteem.

Using multinomial distribution, $E(X) = p_1 \times n$ for each group (denoted as 1 for first group). $Var(X) = n \times p_1(1 - p_1)$ for each group (denoted as 1 for first group).

- 0.45 low self-esteem, expected value: $0.45 \times 850 = 382.5$, variance: $200 \times 0.55(1 - 0.55) = 210.375$
- 0.55 high self-esteem, expected value: $0.55 \times 850 = 467.5$, variance: $200 \times 0.15(1 - 0.15) = 210.375$

Part D

A sample of traffic accidents per day along a 10-mile stretch of I-95 in Virginia that yielded the following results:

Number of Accidents	Frequency
0	121
1	199
2	21
3	12
4	5
5	4
6	2
7	1

Using the poisson distribution, the expected value is equal to the variance: $E(X) = \lambda = Var(X)$, which in this case, from observing the table, is 1.

Question 4

We have been asked to collect 12 signatures for a petition that asks the state government for more money to clean up garbage on public land. The probability of getting a signature from a person approached is 0.40. After finding the mean and variance, answer the following: What is the probability we will have to approach exactly 30 people to get the 12 signatures?

$$P(n) = \binom{n-1}{r-1} p^r (1-p)^{n-r} = P(30) = \binom{30-1}{12-1} (0.40)^{12} (0.60)^{30-12} = 0.0589$$

Question 5

Suppose that we survey six people and find that two of them say they read a newspaper every day and the other four say they do not. We wish to determine the maximum likelihood estimate of p , or the probability of daily newspaper readings among this sample. Use the likelihood function for the binomial distribution to fill in the cell of the following table:

	i $= 2$
p $= 0.1$?
p $= 0.2$?
p $= 0.3$?
p $= 0.4$?

From this table, what is the most likely value of p ?

Using likelihood function for binomial distribution: $P(i) = \binom{n}{i} p^i (1-p)^{n-i}$, $P(i) = \binom{6}{2} p^2 (1-p)^{6-2}$

	$i = 2$
p $= 0.1$	$P(0.1)$ $= 0.0885735$
p $= 0.2$	$P(0.2)$ $= 0.196608$
p $= 0.3$	$P(0.3)$ $= 0.2268945$
p $= 0.4$	$P(0.4)$ $= 0.186624$

The most likely value of p is 0.3.

Question 6

```
In [2]: #!pip install sas7bdat

import numpy as np
import pandas as pd
import statsmodels.api as sm

from patsy import dmatrices
from sas7bdat import SAS7BDAT
```

The Data file USData contains a number of variables from the 50 states in the United States. In this exercise we are interested in using linear regression to predict *violrate*, the rate of violent crimes such as murder, robbery, and assault per 100,000 population in 1995. We shall use the following independent variables: *unemprat* (average monthly unemployment rate in 1995), *density* (population density in 1995), and *gsprod* (gross state product in 1995 -- a measure of the state's economic productivity). Estimate two linear regression models using MLE. The first model is the null model, while the second includes the three independent variables. Use the output from these models to compute the following fit statistics from the second model: McFadden adjusted R^2 and the pseudo- R^2 . Then compute the AIC and BIC from both models.

Loading the Data

```
In [3]: with SAS7BDAT('data/usdata.sas7bdat') as f:
        df_us = f.to_data_frame()
```

```
In [4]: df_us.head(5)
```

Out[4]:

	STATE	ROBBRATE	LARCRATE	ASSRATE	BURGRATE	MURDRATE	FIPS	PERINC	SUICRATE	ASUICRAT	...	UNEMPRAT	POP_
0	Alabama	185.75	2844.27	403.69	1024.83	11.168587	1	19086	12.1	13.2	...	5.1	15223
1	Alaska	155.13	3624.34	526.32	836.92	9.105960	2	NaN	17.1	17.1	...	7.8	24904
2	Arizona	173.76	4925.63	495.71	1416.83	10.407776	4	20068	17.5	18.7	...	5.5	15599
3	Arkansas	125.68	2815.42	379.83	996.90	10.426731	5	17935	14.1	14.5	...	5.4	90281
4	California	331.16	2856.87	590.26	1120.31	11.177942	6	23901	11.1	11.7	...	7.2	11794

5 row s × 25 columns

```
In [5]: y, X = dmatrices('VIOLRATE ~ 1',
                        data=df_us,
                        return_type='dataframe')
model = sm.GLM(y, X)
results = model.fit()
```

```
In [6]: results.summary()
```

Out[6]:

Generalized Linear Model Regression Results

Dep. Variable:	VIOLRATE	No. Observations:	50
Model:	GLM	Df Residuals:	49
Model Family:	Gaussian	Df Model:	0
Link Function:	identity	Scale:	72487.5414926
Method:	IRLS	Log-Likelihood:	-350.22
Date:	Mon, 20 Jun 2016	Deviance:	3.5519e+06
Time:	14:49:42	Pearson chi2:	3.55e+06
No. Iterations:	4		

	coef	std err	z	P> z	[95.0% Conf. Int.]
Intercept	539.4182	38.076	14.167	0.000	464.791 614.045

```
In [7]: results.summary2()
```

Out[7]:

Model:	GLM	AIC:	702.4422
Link Function:	identity	BIC:	3551697.8440
Dependent Variable:	VIOLRATE	Log-Likelihood:	-350.22
Date:	2016-06-20 14:49	LL-Null:	-350.22
No. Observations:	50	Deviance:	3.5519e+06
Df Model:	0	Pearson chi2:	3.55e+06
Df Residuals:	49	Scale:	72488.
Method:	IRLS		

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
Intercept	539.4182	38.0756	14.1670	0.0000	464.7914	614.0450

```
In [8]: y, X = dmatrices('VIOLRATE ~ UNEMPRAT + DENSITY + GSPROD',
                        data=df_us,
                        return_type='dataframe')
model = sm.GLM(y, X)
results = model.fit()
```

```
In [9]: results.summary()
```

Out[9]: Generalized Linear Model Regression Results

Dep. Variable:	VIOLRATE	No. Observations:	50
Model:	GLM	Df Residuals:	46
Model Family:	Gaussian	Df Model:	3
Link Function:	identity	Scale:	50893.5006753
Method:	IRLS	Log-Likelihood:	-339.80
Date:	Mon, 20 Jun 2016	Deviance:	2.3411e+06
Time:	14:49:43	Pearson chi2:	2.34e+06
No. Iterations:	4		

	coef	std err	z	P> z	[95.0% Conf. Int.]
Intercept	62.0352	137.957	0.450	0.653	-208.355 332.425
UNEM PRAT	72.0617	26.258	2.744	0.006	20.596 123.527
DENSITY	0.0421	0.140	0.300	0.764	-0.233 0.317
GSPROD	0.0007	0.000	3.286	0.001	0.000 0.001

```
In [10]: results.summary2()
```

Out[10]:

Model:	GLM	AIC:	687.5993
Link Function:	identity	BIC:	2340921.0780
Dependent Variable:	VIOLRATE	Log-Likelihood:	-339.80
Date:	2016-06-20 14:49	LL-Null:	-350.22
No. Observations:	50	Deviance:	2.3411e+06
Df Model:	3	Pearson chi2:	2.34e+06
Df Residuals:	46	Scale:	50894.
Method:	IRLS		

	Coef.	Std.Err.	z	P> z	[0.025	0.975]
Intercept	62.0352	137.9565	0.4497	0.6529	-208.3546	332.4251
UNEM PRAT	72.0617	26.2584	2.7443	0.0061	20.5962	123.5272
DENSITY	0.0421	0.1404	0.3001	0.7641	-0.2330	0.3172
GSPROD	0.0007	0.0002	3.2864	0.0010	0.0003	0.0011