

# Homework 2

Jacob Puthipiroj

March 18, 2020

## 1 A Theory of Extramarital Affairs

- (a) The regressors of interest are  $v_1$  to  $v_8$ ; however, not necessarily all of them belong in your model. Use these data to build a binary choice model for  $A$ . Report all computed results for the model. Compute the marginal effects for the variables you choose. Compare the results you obtain for a probit model to those for a logit model. Are there any substantial differences in the results for the two models?

The specification of the probit model is

$$p = \Phi(\beta_0 + \mathbf{B}x), \quad \Phi(x) = \int_{-\infty}^x \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2} dx$$

The specification of the logit model is

$$p = \frac{1}{1 + e^{-l}} = \frac{e^l}{1 + e^l}, \quad l = \ln\left(\frac{p}{1-p}\right) = \beta_0 + \mathbf{B}x$$

In both cases,  $\beta_0$  is the intercept, and we define  $\mathbf{B}x$  as

$$\mathbf{B}x = \sum_{i=1}^6 \beta_i v_i + \sum_{i=1}^6 \delta_{i,occ} v_{i,occ} + \sum_{i=1}^6 \delta_{i,hoc} v_{i,hoc}$$

Thus all variables  $v_1$  through  $v_6$  (rating, age, years, children, religiosity, and education) are thought to be continuous variables, while the occupation of the wife and husband are considered indicator variables, with each  $\delta_i$  corresponding to each of the 6 occupation categories. The output of the probit regression is available here, and the output of the logit regression here in the appendix. In both the probit and logit model, the same coefficients were considered significant at the  $\alpha = 0.05$  level: rating, age, years, religiosity, and the two dummy variables for the occupation of the wife being in a managerial, administrative or business role  $\delta_{5,occ}$  and being a professional with an advanced degree,  $\delta_{6,occ}$ .

The negative signs for  $\beta_{\text{rating}}$  and  $\beta_{\text{religiosity}}$  make sense: The better one rates [satisfaction of] a marriage, the less reason there is to be involved in an extramarital affair. Similarly, the more religious one is, the less one is less likely to cheat, perhaps due to religious beliefs concerning the (im)morality of infidelity. The negative sign for  $\beta_{\text{age}}$  implies that older women cheat less, when controlled for the years of the marriage (both are highly correlated with each other at  $r = 0.8941$ ).  $\beta_{\text{years}}$  is positive, perhaps because as a marriage drags on, one is more inclined to look elsewhere for emotional or physical fulfillment, such as an extramarital affair. A final interesting note is that females in managerial and professional careers are more likely to have cheated if they were a student.

The marginal effects of each variable vary at different points, and as such we compute instead the average marginal effects (AME) by calculating the marginal effect for each individual with their observed levels of covariates, which are then averaged across individuals. The AMEs for the variables in the probit model is available here in the appendix. The marginal effects of the logit model were similar, available here in the appendix. There were no substantial differences in the results between the two models. On average, each additional increase in the reported marriage score decreased the probability of cheating by around 13%, each additional year of the female decreased the probability by 1.1%, each additional year of marriage increased the probability by 2%, each additional score of religiosity decreased the probability by 6.8%, being in a managerial, administrative or business occupation increased the probability of extramarital affairs by 18%, and being in a professional career with an advanced degree increased the probability by 19%, both compared to the baseline of being a student.

- (b) Continuing the analysis from part a), we now consider the self-reported rating,  $v1$ . This is a natural candidate for an ordered choice model, because the simple five-item coding is a censored version of what would be a continuous scale on some subjective satisfaction variable. Analyze this variable using an ordered probit model. What variables appear to explain the response to this survey question? Can you obtain the marginal effects for your model? Report them as well. What do they suggest about the impact of the different independent variables on the reported ratings?

The output for the ordered probit regression is available here. As seen, the only significant predictors (at the  $\alpha = 0.05$  level) were children, religiosity, and the husband's occupation being in farming, agriculture, semi-skilled, or unskilled labor.  $\beta_{\text{children}}$  was negative, suggesting that more children made marriages unhappier, while  $\beta_{\text{religiosity}} > 0$  suggests that more religiosity made marriages happier.  $\delta_{2,hoc} < 0$ , suggesting that wives were less satisfied with their marriage when their husbands were involved in menial labor, compared to when they (the husbands) were students.

The average marginal effects are available here. An interesting finding suggested that for people whose marriages were given a rating of 1, 2, 3 or 4, each additional child was found to increase marriage satisfaction by 0.2%, 0.6%, 1%, and 0.5% respectively. However, when the marriage was rated 5, each additional child was found to decrease marriage satisfaction by 2.4%, creating a negative feedback loop and downward pressure away from the maximum satisfaction score. The opposite pattern was found for religiosity: each additional unit of religiosity decreased marriage satisfaction by 0.5%, 1.2%, 2%, and 1.1% for marriages rated 1, 2, 3 or 4 respectively, but increased satisfaction by 5% for marriages already rated at 5, creating a positive feedback loop. The average marginal effects for husbands having a menial labor occupation followed a similar pattern as that of children: 0.6%, 1.5% and 2.8% for marriages rated 1, 2 or 3 respectively, not significant for marriages rated 4, and  $-6.6\%$  for marriages rated 5.

This suggests that different independent variables have different marginal effects for the dependent variable being at different levels. Having additional children, or having the husband pick up a menial labor job, for example, acts as a normalizer, and could improve the satisfaction of a failing marriage, but prevents the marriage from achieving the highest satisfaction level. On the other hand, religiosity would seem to be a polarizer, making bad marriages worse, and great marriages better.

## 2 Incentive Effects in the Demand for Health Care

A note about this dataset: there were several mistakes in some of the data. The indicator variable handdum for example, was miscoded in the year 1987. Ones and zeros were swapped, and if not corrected, would imply that 88% of the participants in 1987 were handicapped, when in reality it was  $100 - 88 = 12\%$ . The full list of coding errors are given here.[1]

- (a) Begin by fitting a Poisson model to this variable. The exogenous variables are listed in Table F7.1. Determine an appropriate specification for the right-hand side of your model. Report the regression results and the marginal effects.

The specification we have is

$$\mathbf{B}\mathbf{x} = \sum_{i=1}^4 \beta_i v_i + \sum_{i=1}^3 \delta_i u_i$$

Where  $v_i$  indicate each of the four continuous variables of age, hsat, educ and docvis, and  $u_i$  indicate each of the dummy variables handdum, addon, and bluec. The Poisson regression results are available here in the appendix. All variables, save for the constant, were found to be significant. Specifically, the coefficients for age, hsat, educ and bluec were negative, while the coefficients for docvis, handdum, and addon were positive.

The positive signs for  $\beta_{\text{docvis}}$ ,  $\delta_{\text{handdum}}$  and  $\delta_{\text{addon}}$  are unsurprising. Individuals who have frequent doctor visits in the past 3 months are likely to also frequently visit hospitals in the last calendar year. Handicapped individuals may need to visit hospitals more, as a result of a chronic treatment for said handicap, or a new injury causing patients to become handicapped and require treatment. Finally, those who purchase add-on insurance may expect themselves to be more at risk of requiring healthcare.

The negative signs for  $\beta_{\text{hsat}}$  and  $\beta_{\text{educ}}$  also make sense. Individuals who are more satisfied with their health are likely healthier, and require fewer hospital visits. More educated individuals may also make wiser health-related decisions, such as in career and lifestyle choices, and visit hospitals less often. The negative sign for  $\delta_{\text{bluec}}$ , however, is surprising, as workplace injuries (which would require hospital visits) should be more commonplace among blue collar workers. Several explanations are possible: Perhaps only physically healthier people would consider blue collar jobs, or, blue collar workers choose to visit hospitals only for the most serious injuries, either out of necessity or because of being desensitized to workplace injuries.  $\beta_{\text{age}} < 0$  is similarly surprising. Despite (or because of) having weaker bodies, older people may make lifestyle and professional choices that would limit physical injuries.

The marginal effects are given in the appendix here. Again, the marginal effects of all explanatory variables were significant. On average, being handicapped and being insured by addon insurance increased the expected number of hospital visits by 4% and 5% respectively, and each additional doctor visit in the past 3 months increased expected visits by 0.4%. Each additional year in age and in schooling decreased expected visits by 0.09% and 0.6% respectively. Each additional unit increase of perceived health satisfaction on the 1-10 scale decreased expected visits by 2.5%. Finally, being in a blue-collar job decreased expected visits by 1.4%.

- (b) Estimate the model using ordinary least squares and compare your least squares results to the marginal effects computed in part a). What do you find?

The output is given here in the appendix. Unlike in the Poisson regression, the variables `addon` and `bluec` were found to be insignificant. In the OLS model with no interaction terms, the coefficient estimates are equivalent to the marginal effects, and can be compared with the AME of the Poisson regression.

Like in the Poisson model, being handicapped had large marginal effects on increasing the expected number of visits, at  $\hat{\delta}_{\text{handdum}} = 5.8\%$ , comparable to the 4% from the Poisson model. However, the coefficients for being on `addon` insurance, and in a blue-collar occupation were not significant. Each additional doctor visit in the past 3 months increased expected visits by 1.7%, comparable to the 4% from the Poisson model. The marginal effects for each additional year in age and schooling was  $-0.15\%$  and  $-0.5\%$  respectively, compared to  $-0.09\%$  and  $0.6\%$  from the Poisson model. Each additional unit increase of perceived health satisfaction on the 1-10 scale decreased expected visits by 2.4%, similar to 2.5% from the Poisson model.

Overall, the marginal effect estimates were identical in sign and very similar in magnitude, though  $\delta_{\text{addon}}$  and  $\delta_{\text{bluec}}$  were not significant like they were in the Poisson model.

- (c) Is there evidence of overdispersion in the data? Test for overdispersion.

Overdispersion is where variance is greater than would be expected in a Poisson regression. One test for this is to simply run the Poisson regression, and then test using the `poisgof` command. The output is

```
Deviance goodness-of-fit = 20004.18
Prob > chi2(27322)       = 1.0000

Pearson goodness-of-fit   = 131234
Prob > chi2(27322)       = 0.0000
```

Which suggests the Poisson model to be inappropriate. We can furthermore check for overdispersion by using the `nbreg` command, given here in the appendix, which fits the data with a negative binomial distribution, and gives a likelihood-ratio test, with the hypothesis being that the negative binomial distribution is equivalent to a poisson distribution, under  $\alpha = 0$ . Since the p-value is significant,  $\alpha$  is significantly different from 0, so the poisson model is again inappropriate.

### 3 Appendix

The accompanying code is given as a Jupyter Notebook (in Stata) as well as a .do file.

#### Probit Model Output

Probit regression	Number of obs	=	6,366
	LR chi2(16)	=	1097.04
	Prob > chi2	=	0.0000
Log likelihood = -3454.0116	Pseudo R2	=	0.1370

A	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
rating	-.42506	.0183553	-23.16	0.000	-.4610357	-.3890843
age	-.0359131	.0060622	-5.92	0.000	-.0477948	-.0240314
years	.0642991	.006468	9.94	0.000	.0516221	.0769761
children	.0086311	.0190523	0.45	0.651	-.0287107	.045973
religiosity	-.2235751	.0205278	-10.89	0.000	-.2638088	-.1833414
education	-.001846	.0102593	-0.18	0.857	-.0219539	.0182618
occupation						
2	.2107086	.2498365	0.84	0.399	-.2789619	.7003792
3	.3992348	.2461618	1.62	0.105	-.0832334	.881703
4	.2630816	.2467651	1.07	0.286	-.2205691	.7467324
5	.6091108	.2495332	2.44	0.015	.1200347	1.098187
6	.6480159	.2792565	2.32	0.020	.1006832	1.195349
husbandocc						
2	.0943557	.1069709	0.88	0.378	-.1153034	.3040147
3	.1662867	.1171215	1.42	0.156	-.0632672	.3958406
4	.0766486	.1039138	0.74	0.461	-.1270187	.280316
5	.0915583	.1049421	0.87	0.383	-.1141246	.2972411
6	.1029664	.1173516	0.88	0.380	-.1270385	.3329713
_cons	1.792191	.325053	5.51	0.000	1.155099	2.429283

Average marginal effects                      Number of obs        =        6,366  
Model VCE        : OIM

	Delta-method					
	dy/dx	Std. Err.	z	P> z	[95% Conf. Interval]	
rating	-.1302394	.0049252	-26.44	0.000	-.1398926	-.1205863
age	-.0110039	.0018444	-5.97	0.000	-.0146188	-.0073889
years	.0197014	.0019401	10.15	0.000	.0158989	.0235039
children	.0026446	.0058375	0.45	0.651	-.0087967	.0140859
religiosity	-.068504	.0061325	-11.17	0.000	-.0805234	-.0564845
education	-.0005656	.0031434	-0.18	0.857	-.0067266	.0055954
occupation						
2	.056928	.0634788	0.90	0.370	-.0674882	.1813441
3	.1136303	.0624022	1.82	0.069	-.0086757	.2359363
4	.072174	.06255	1.15	0.249	-.0504218	.1947697
5	.1818514	.0639968	2.84	0.004	.0564201	.3072828
6	.1949477	.0768603	2.54	0.011	.0443043	.345591
husbandocc						
2	.0283895	.0317009	0.90	0.370	-.0337431	.0905221
3	.0508015	.0352348	1.44	0.149	-.0182575	.1198605
4	.0229711	.0307046	0.75	0.454	-.0372087	.083151
5	.0275308	.0310522	0.89	0.375	-.0333304	.088392
6	.0310389	.0350368	0.89	0.376	-.037632	.0997098

## Logit Model Output

Logistic regression

Number of obs = 6,366

LR chi2(16) = 1092.71

Prob > chi2 = 0.0000

Pseudo R2 = 0.1365

Log likelihood = -3456.1733

A	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
rating	-.7102283	.0314818	-22.56	0.000	-.7719314	-.6485252
age	-.0612789	.0103231	-5.94	0.000	-.0815117	-.041046
years	.107976	.0109772	9.84	0.000	.0864611	.1294909
children	.0156448	.0320509	0.49	0.625	-.0471737	.0784634
religiosity	-.3753863	.0348686	-10.77	0.000	-.4437274	-.3070451
education	-.0017253	.017398	-0.10	0.921	-.0358247	.032374
occupation						
2	.3902386	.4475507	0.87	0.383	-.4869446	1.267422
3	.7026792	.4414598	1.59	0.111	-.1625661	1.567925
4	.4713969	.4425232	1.07	0.287	-.3959326	1.338726
5	1.054197	.4466347	2.36	0.018	.1788096	1.929585
6	1.108015	.4942125	2.24	0.025	.1393766	2.076654
husbandocc						
2	.170447	.1860868	0.92	0.360	-.1942766	.5351705
3	.2841727	.2021606	1.41	0.160	-.1120548	.6804002
4	.1428406	.1810041	0.79	0.430	-.211921	.4976022
5	.1723288	.1826377	0.94	0.345	-.1856345	.5302921
6	.1827633	.203652	0.90	0.369	-.2163872	.5819138
_cons	2.970755	.5722101	5.19	0.000	1.849244	4.092266

Average marginal effects	Number of obs	=	6,366
Model VCE : OIM			

	Delta-method					
	dy/dx	Std. Err.	z	P> z	[95% Conf. Interval]	
rating	-.1290391	.0048629	-26.54	0.000	-.1385703	-.1195079
age	-.0111336	.0018588	-5.99	0.000	-.0147768	-.0074904
years	.0196178	.001944	10.09	0.000	.0158077	.023428
children	.0028425	.005823	0.49	0.625	-.0085704	.0142553
religiosity	-.0682027	.0061495	-11.09	0.000	-.0802556	-.0561499
education	-.0003135	.003161	-0.10	0.921	-.0065089	.0058819
occupation						
2	.0612388	.065116	0.94	0.347	-.0663862	.1888639
3	.1168366	.0640436	1.82	0.068	-.0086865	.2423598
4	.0751529	.0642053	1.17	0.242	-.0506871	.200993
5	.185137	.0655919	2.82	0.005	.0565792	.3136947
6	.1960177	.0785268	2.50	0.013	.0421081	.3499274
husbandocc						
2	.0302712	.0324327	0.93	0.351	-.0332957	.0938381
3	.0513025	.0358087	1.43	0.152	-.0188812	.1214862
4	.0252634	.0314419	0.80	0.422	-.0363616	.0868885
5	.030614	.0317795	0.96	0.335	-.0316727	.0929006
6	.0325179	.0358017	0.91	0.364	-.0376522	.102688



## Ordered Probit Model

Ordered probit regression	Number of obs	=	6,366
	LR chi2(15)	=	236.49
	Prob > chi2	=	0.0000
Log likelihood = -7808.2421	Pseudo R2	=	0.0149

rating	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	-.0047552	.0047154	-1.01	0.313	-.0139971	.0044867
years	-.0070395	.0050613	-1.39	0.164	-.0169594	.0028804
children	-.0632364	.0153484	-4.12	0.000	-.0933187	-.0331541
religiosity	.1310093	.0161123	8.13	0.000	.0994298	.1625888
education	.0140007	.0081484	1.72	0.086	-.0019699	.0299713
occupation						
2	-.1317628	.1824797	-0.72	0.470	-.4894165	.2258909
3	-.2022679	.1794037	-1.13	0.260	-.5538926	.1493568
4	-.0632041	.1798682	-0.35	0.725	-.4157393	.2893311
5	-.1434889	.1826641	-0.79	0.432	-.5015038	.2145261
6	-.2022114	.2097835	-0.96	0.335	-.6133796	.2089568
husbandocc						
2	-.1708474	.0822583	-2.08	0.038	-.3320708	-.0096241
3	-.1705628	.0907884	-1.88	0.060	-.3485048	.0073792
4	-.0981205	.0797453	-1.23	0.219	-.2544184	.0581774
5	-.0671601	.0806878	-0.83	0.405	-.2253052	.090985
6	.0409246	.0912108	0.45	0.654	-.1378452	.2196944
/cut1	-2.221501	.2388776			-2.689692	-1.753309
/cut2	-1.522325	.2363926			-1.985646	-1.059004
/cut3	-.7805359	.2356443			-1.24239	-.3186815
/cut4	.1906019	.2355468			-.2710615	.6522652



## Poisson Model

Poisson regression

Number of obs = 27,308

LR chi2(7) = 2123.85

Prob > chi2 = 0.0000

Pseudo R2 = 0.0791

Log likelihood = -12362.865

hospvis	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	-.0071362	.0015424	-4.63	0.000	-.0101592	-.0041132
hsat	-.186052	.0069948	-26.60	0.000	-.1997616	-.1723424
educ	-.045944	.0084623	-5.43	0.000	-.0625299	-.0293582
docvis	.0306762	.0011552	26.55	0.000	.028412	.0329404
1.handdum	.2691112	.0439668	6.12	0.000	.1829378	.3552846
1.addon	.3218622	.1077161	2.99	0.003	.1107425	.5329819
1.bluec	-.1103086	.0405875	-2.72	0.007	-.1898586	-.0307585
_cons	-.1971156	.1345566	-1.46	0.143	-.4608416	.0666105

Average marginal effects	Number of obs	=	27,308
Model VCE : OIM			

	Delta-method					
	dy/dx	Std. Err.	z	P> z	[95% Conf. Interval]	
age	-.0009865	.0002138	-4.61	0.000	-.0014056	-.0005674
hsat	-.0257194	.0010537	-24.41	0.000	-.0277846	-.0236543
educ	-.0063512	.0011744	-5.41	0.000	-.0086529	-.0040495
docvis	.0042406	.000174	24.38	0.000	.0038996	.0045816
1.handdum	.0403667	.0071613	5.64	0.000	.0263309	.0544026
1.addon	.0521476	.0202966	2.57	0.010	.012367	.0919282
1.bluec	-.0148143	.0052965	-2.80	0.005	-.0251954	-.0044333

## OLS Regression

Source	SS	df	MS	Number of obs	=	27,308
-----+-----				F(7, 27300)	=	91.40
Model	489.306593	7	69.9009419	Prob > F	=	0.0000
Residual	20877.8454	27,300	.764756242	R-squared	=	0.0229
-----+-----				Adj R-squared	=	0.0226
Total	21367.152	27,307	.782478925	Root MSE	=	.8745

-----+-----						
hospvis	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]	
-----+-----						
age	-.0015269	.0005016	-3.04	0.002	-.0025101	-.0005436
hsat	-.0240788	.0026008	-9.26	0.000	-.0291765	-.0189811
educ	-.0049511	.0024185	-2.05	0.041	-.0096916	-.0002106
docvis	.0170985	.0010107	16.92	0.000	.0151175	.0190795
1.handdum	.0580785	.0183262	3.17	0.002	.0221582	.0939988
1.addon	.0443904	.0390841	1.14	0.256	-.0322164	.1209971
1.bluec	-.0159247	.0128963	-1.23	0.217	-.0412021	.0093527
_cons	.3663938	.0439651	8.33	0.000	.28022	.4525676
-----+-----						

## References

- [1] Regina T Riphahn, Achim Wambach, and Andreas Million. Incentive effects in the demand for health care: a bivariate panel count data estimation. *Journal of applied econometrics*, 18(4):387–405, 2003.

## Negative Binomial Model

Negative binomial regression	Number of obs	=	27,326
	LR chi2(3)	=	676.70
Dispersion = mean	Prob > chi2	=	0.0000
Log likelihood = -10037.84	Pseudo R2	=	0.0326

hospvis	Coef.	Std. Err.	z	P> z	[95% Conf. Interval]	
age	-.006069	.0022079	-2.75	0.006	-.0103964	-.0017417
hsat	-.2193225	.0099992	-21.93	0.000	-.2389206	-.1997244
handdum	.4574462	.0707898	6.46	0.000	.3187007	.5961916
_cons	-.4580444	.1258712	-3.64	0.000	-.7047473	-.2113414
/lnalpha	1.896108	.0403598			1.817004	1.975212
alpha	6.659923	.2687934			6.153396	7.208146

LR test of alpha=0: chibar2(01) = 5175.62      Prob >= chibar2 = 0.000

## References

- [1] Regina T Riphahn, Achim Wambach, and Andreas Million. Incentive effects in the demand for health care: a bivariate panel count data estimation. *Journal of applied econometrics*, 18(4):387–405, 2003.