

# Linear Probability Model and Logistic Regression

Andy Grogan-Kaylor

28 Dec 2020 08:53:39

## Introduction

The *Linear Probability Model* (LPM) is often discussed as an alternative to *logistic* regression. Essentially, the LPM is a linear model with a *dichotomous* dependent variable.

## Setup

```
. clear all

. use http://www.stata-press.com/data/r15/margex, clear // artificial data from Stata
(Artificial data for margins)
```

## Background

I read through a number of references to develop this handout, especially the excellent book on Categorical Data Analysis by Long and Freese, and the always excellent Stata documentation. As I was finishing up this handout, I came across a superb handout by Richard Williams (referenced below), which does a better and more thorough job of explaining these issues than this short handout. You are encouraged to look it up.

*Broadly speaking* the Linear Probability Model is likely to give similar results to the logistic regression model:  $\beta$  coefficients are likely to have the same directions and similar statistical significances.

However, as one compares these approaches *more closely*, the Linear Probability Model is arguably incorrect on several grounds, some of which are illustrated in the figure below:

```
. twoway (lowess outcome age) (lfit outcome age), ///
> title("Outcome By Age") ///
> legend(order(1 "lowess smoother" 2 "linear fit")) ///
> scheme(michigan)

. graph export mygraph0.png, width(1000) replace
(file /Users/agrogan/Desktop/newstuff/categorical/LPM-and-logistic/mygraph0.png written in PN
> G format)
```

1. Marginal effects are mis-stated: The smoother indicates that the relationship of outcome and age is curvilinear. Thus, the effect on  $y$  of a 1 unit increase in  $x$  is different for different values of  $x$ .
2. Predictions can be implausible: By definition, negative probabilities are clearly impossible. However the linear fit predicts negative probabilities of the outcome for lower values of age.

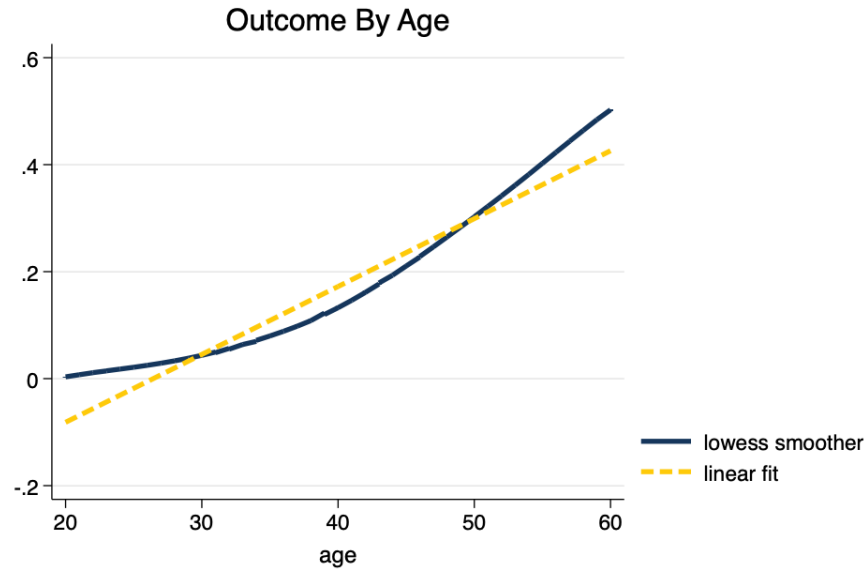


Figure 1: Lowess Smoother and Linear Fit Of Outcome By Age

3. Data with a dichotomous outcome are by definition heteroskedastic. The LPM (unless corrections are applied) makes assumptions of homoskedasticity. Thus, inferences about statistical significance—or the lack thereof—are likely to be incorrect.

These differences in results are likely to become more salient the more one pays *detailed attention* to marginal effects for different values of the independent variables, and to predicted probabilities for different values of the independent variables.

## Compare LPM and Logistic Regression In More Detail

### Confirm That Outcome Is Dichotomous

```
. tabulate outcome // outcome is dichotomous
```

outcome	Freq.	Percent	Cum.
0	2,491	83.03	83.03
1	509	16.97	100.00
Total	3,000	100.00	

### Linear Probability Model

```
. regress outcome sex##c.age i.group // linear probability model
```

Source	SS	df	MS	Number of obs	=	3,000
Model	79.386424	5	15.8772848	F(5, 2994)	=	138.49
Residual	343.253243	2,994	.114647042	Prob > F	=	0.0000
				R-squared	=	0.1878
				Adj R-squared	=	0.1865
Total	422.639667	2,999	.140926865	Root MSE	=	.3386

outcome	Coef.	Std. Err.	t	P> t	[95% Conf. Interval]
---------	-------	-----------	---	------	----------------------

sex						
female	-.2320346	.0489015	-4.74	0.000	-.3279185	-.1361508
age	.0061307	.0008814	6.96	0.000	.0044025	.0078589
sex#c.age						
female	.0072707	.0011613	6.26	0.000	.0049936	.0095477
group						
2	-.0888273	.0164698	-5.39	0.000	-.1211206	-.0565339
3	-.1034404	.0220694	-4.69	0.000	-.1467131	-.0601676
_cons	-.0597978	.0401266	-1.49	0.136	-.1384763	.0188806

```
. predict yhat_LPM // predicted probabilities
(option xb assumed; fitted values)

. twoway scatter yhat_LPM age, ///
> title("Predicted Probabilities from Linear Probability Model") ///
> scheme(michigan)

. graph export myLPM.png, width(1000) replace
(file /Users/agrogan/Desktop/newstuff/categorical/LPM-and-logistic/myLPM.png written in PNG f
> ormat)
```

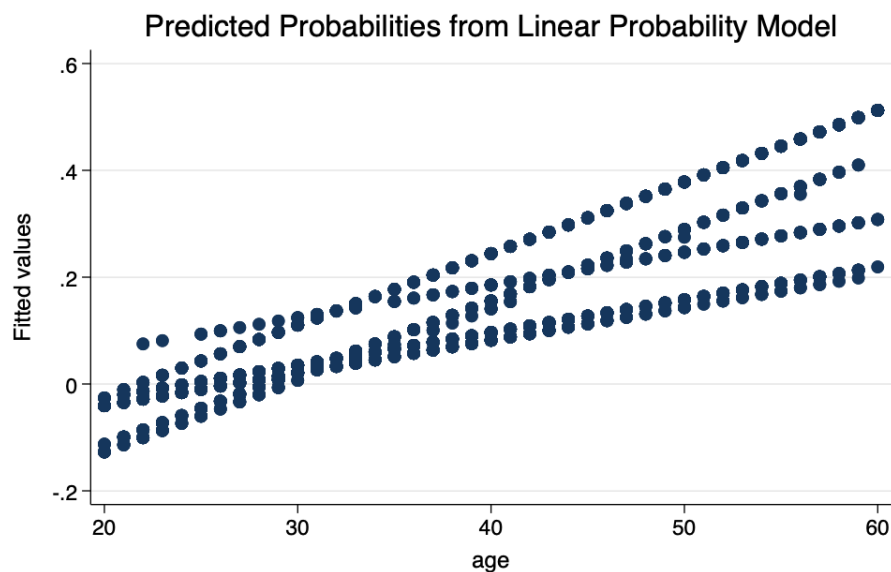


Figure 2: Predicted Values from Linear Probability Model

## Logistic Regression

```
. logit outcome sex#c.age i.group // logistic regression model
Iteration 0: log likelihood = -1366.0718
Iteration 1: log likelihood = -1118.129
Iteration 2: log likelihood = -1070.8227
Iteration 3: log likelihood = -1068.0102
Iteration 4: log likelihood = -1067.99
Iteration 5: log likelihood = -1067.99

Logistic regression      Number of obs      =      3,000
                        LR chi2(5)      =      596.16
```



Williams, R. (2015). *Logistic Regression, Part I: Problems with the Linear Probability Model (LPM)*. South Bend, IN.