



LECTURE 5: TRICKY DVs

Dr. Rachel Blum

November 12, 2019

BEYOND OLS

MOTIVATION: INFERENCE

- We have theories. We want to know how “true” they are.
- We have data. We want to know how that data could help us evaluate our theories.

MOTIVATION: INFERENCE

- We have theories. We want to know how “true” they are.
- We have data. We want to know how that data could help us evaluate our theories.
- One key approach to these questions comes from **Bayes’ Theorem**:

$$P(\text{hypothesis}|\text{data}) = \frac{P(\text{data}|\text{hypothesis}) \times P(\text{hypothesis})}{P(\text{data})}$$

$$P(\text{hypothesis}|\text{data}) = \frac{P(\text{data}|\text{hypothesis}) \times P(\text{hypothesis})}{P(\text{data})}$$

- The problem is, we don't know $P(\text{hypothesis})$.
- $P(\text{data}|\text{hypothesis})$ is not so difficult. Recall the p-value.
- $P(\text{data})$ is calculable. It's a function of $P(\text{data}|\text{hypothesis})$ and $P(\text{hypothesis})$ and it drops out.

BAYES THEOREM EXAMPLE

Suppose we want to catch terrorists. We know that a large percentage of terrorists are Muslim. Should we not target Muslims in attempting to catch terrorists?

BAYES THEOREM EXAMPLE

BAYES THEOREM EXAMPLE

- Probability that a terrorist is a Muslim: 95%
 - A lot depends on how we define 'Muslim' and especially 'terrorist.'

BAYES THEOREM EXAMPLE

- Probability that a terrorist is a Muslim: 95%
 - A lot depends on how we define 'Muslim' and especially 'terrorist.'
- Probability that a person is a Muslim: 23%
 - Smaller in United States or in western airports, etc.

BAYES THEOREM EXAMPLE

- Probability that a terrorist is a Muslim: 95%
 - A lot depends on how we define 'Muslim' and especially 'terrorist.'
- Probability that a person is a Muslim: 23%
 - Smaller in United States or in western airports, etc.
- Probability that a person is a terrorist: 0.01%
 - An absurdly large "guess".

BAYES THEOREM EXAMPLE

$P(\text{Muslim}) = \text{probability of being Muslim} = 0.23$

$P(\text{terrorist}) = \text{probability of being a terrorist} = 0.0001$

$P(\text{Muslim}|\text{terrorist}) = \text{probability a terrorist is Muslim} = 0.95$

BAYES THEOREM EXAMPLE

$P(\text{Muslim}) = \text{probability of being Muslim} = 0.23$

$P(\text{terrorist}) = \text{probability of being a terrorist} = 0.0001$

$P(\text{Muslim}|\text{terrorist}) = \text{probability a terrorist is Muslim} = 0.95$

We want to know $P(\text{terrorist}|\text{Muslim})$

$$\begin{aligned} P(\text{terrorist}|\text{Muslim}) &= \frac{P(\text{Muslim}|\text{terrorist}) \times P(\text{terrorist})}{P(\text{Muslim})} \\ &= \frac{(0.95)(0.0001)}{(0.23)} \\ &= \frac{0.114}{0.23} = 0.0004 \end{aligned}$$

BAYES THEOREM EXAMPLE

$P(\text{Muslim}) = \text{probability of being Muslim} = 0.23$

$P(\text{terrorist}) = \text{probability of being a terrorist} = 0.0001$

$P(\text{Muslim}|\text{terrorist}) = \text{probability a terrorist is Muslim} = 0.95$

We want to know $P(\text{terrorist}|\text{Muslim})$

$$\begin{aligned}P(\text{terrorist}|\text{Muslim}) &= \frac{P(\text{Muslim}|\text{terrorist}) \times P(\text{terrorist})}{P(\text{Muslim})} \\ &= \frac{(0.95)(0.0001)}{(0.23)} \\ &= \frac{0.114}{0.23} = 0.0004\end{aligned}$$

or a 0.004% chance that a Muslim is a terrorist.

- **Bayesian statistics** attempts to use your beliefs about $P(\text{hypothesis})$ as a basis for calculating $P(\text{hypothesis}|\text{data})$.

- **Bayesian statistics** attempts to use your beliefs about $P(\text{hypothesis})$ as a basis for calculating $P(\text{hypothesis}|\text{data})$.
 - Maximum Likelihood framework says $P(\text{data}|\text{hypothesis})$ is at least proportional to $P(\text{hypothesis}|\text{data})$, so $P(\text{data}|\text{hypothesis})$ is all we really care about.

BAYESIAN AND FREQUENTISM COMPARED

- **Bayesian statistics** attempts to use your beliefs about $P(\text{hypothesis})$ as a basis for calculating $P(\text{hypothesis}|\text{data})$.
 - Maximum Likelihood framework says $P(\text{data}|\text{hypothesis})$ is at least proportional to $P(\text{hypothesis}|\text{data})$, so $P(\text{data}|\text{hypothesis})$ is all we really care about.
- Frequentist statistics (OLS) will just learn to live with $P(\text{data}|\text{hypothesis})$

BAYESIAN AND FREQUENTISM COMPARED

- **Bayesian statistics** attempts to use your beliefs about $P(\text{hypothesis})$ as a basis for calculating $P(\text{hypothesis}|\text{data})$.
 - Maximum Likelihood framework says $P(\text{data}|\text{hypothesis})$ is at least proportional to $P(\text{hypothesis}|\text{data})$, so $P(\text{data}|\text{hypothesis})$ is all we really care about.
- Frequentist statistics (OLS) will just learn to live with $P(\text{data}|\text{hypothesis})$
 - Hypothesis testing: Reject null hypothesis when $P(\text{data}|\text{null hypothesis})$ is small.

WHEN OLS DOESN'T WORK

WHEN OLS WORKS (GAUSS MARKOFF)

When correctly specified, OLS is the **Best Linear Unbiased Estimator** of the regression line.

The model is correctly specified.

$$Y = \beta_0 + \beta_1 X_1 + \epsilon$$

Errors are distributed with mean zero.

$$E(\epsilon) = 0$$

Errors have constant variance (homoskedasticity).

$$\text{Var}(\epsilon) = \sigma^2 (\text{vs. } \sigma_i^2)$$

Errors are uncorrelated with each other (no autocorrelation).

$$\text{Cov}(\epsilon_i, \epsilon_j) = 0, i \neq j$$

Errors are uncorrelated with the independent variables.

$$\text{Cov}(\epsilon_i, X_i) = 0$$

Helpful: errors are normally distributed, especially for small N).

$$\epsilon \sim N(0, \sigma_2)$$

The dependent variable is a linearly additive function of a set of linearly independent independent variables, plus a disturbance term.

There is no error associated with the measurement of the IV (X).

Regressors are linearly independent (i.e., *not* multicollinear).

You have more observations than regressors. If

$$X = n \times k \implies n > k$$

OLS only works if you are modeling a linear relationship.

CONSEQUENCES OF GAUSS MARKOFF VIOLATIONS

If linearity assumptions are violated \Rightarrow we have estimated the wrong thing.

Results are meaningless.

CONSEQUENCES OF GAUSS MARKOFF VIOLATIONS

If autocorrelation of errors occurs \Rightarrow standard errors are wrong.

Sign of mis-specification.

CONSEQUENCES OF GAUSS MARKOFF VIOLATIONS

If key variable is omitted \Rightarrow estimates will be incorrect.

Standard errors (and p -values, etc.) will also be incorrect.

NON-LINEARITY: DICHOTOMOUS OUTCOMES

DICHOTOMOUS OUTCOMES

- A **dichotomous outcome** is divided in two parts: the outcome either occurred or it did not.

DICHOTOMOUS OUTCOMES

- A **dichotomous outcome** is divided in two parts: the outcome either occurred or it did not.
 - A policy intervention either works or does not work.

DICHOTOMOUS OUTCOMES

- A **dichotomous outcome** is divided in two parts: the outcome either occurred or it did not.
 - A policy intervention either works or does not work.
 - Votes can be either in favor of or against a policy or representative.

DICHOTOMOUS OUTCOMES

- A **dichotomous outcome** is divided in two parts: the outcome either occurred or it did not.
 - A policy intervention either works or does not work.
 - Votes can be either in favor of or against a policy or representative.
 - Survey respondents are either for or against a given prompt

ENTER: MAXIMUM LIKELIHOOD ESTIMATION (MLE)

- **Practical motivation:** A model for when OLS might not work so well.
- **Theoretical motivation:** Making inferences based on probabilities.
- **The mechanics:** What are we maximizing the likelihood of, and how?

MLE is a different estimator for the parameters of a model (another way to estimate β .)

Especially helpful when the relationship you're estimating is not linear.

- Both **probit** and **logit** models estimate & predict probable relationship between Y and X .

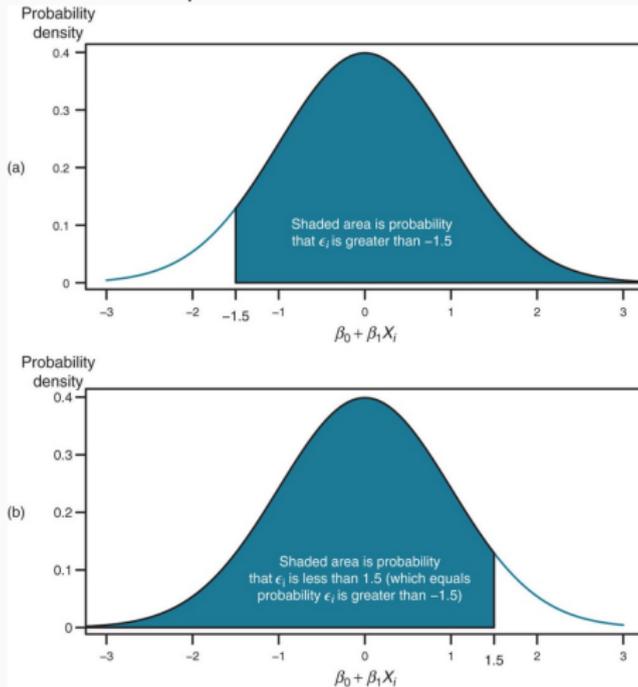
- Both **probit** and **logit** models estimate & predict probable relationship between Y and X .
- Predicted probabilities range between 0 and 1.

- Both **probit** and **logit** models estimate & predict probable relationship between Y and X .
- Predicted probabilities range between 0 and 1.
- Achieved by fitting **S** curve to the data

- Both **probit** and **logit** models estimate & predict probable relationship between Y and X .
- Predicted probabilities range between 0 and 1.
- Achieved by fitting **S** curve to the data
- Practically the same (but we will consider each in turn)

TWO MODELS FOR DICHOTOMOUS DVs: LOGIT AND PROBIT

Logit and Probit both provide a curve **bounded** between 0 and 1.



- **Probit** model is one way to estimate the effects of IVs on a dichotomous outcome.
- **Key assumption:** error term ϵ has a **normal distribution**.
- Model estimates the **probability** of Y taking a certain value, given whatever X 's value is:

$$Pr(Y = 1|X_1) = Pr(\epsilon \leq \beta_0 + \beta_1 X_1)$$

NORMAL DISTRIBUTION AND CDF

CDF = *Cumulative Density Function* of a normal distribution.

NORMAL DISTRIBUTION AND CDF

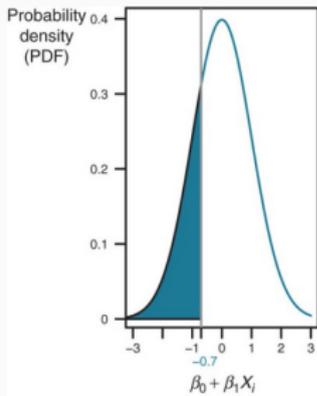
CDF = *Cumulative Density Function* of a normal distribution.

CDF written as probability of something being $<$ a certain value in a *normal distribution*.

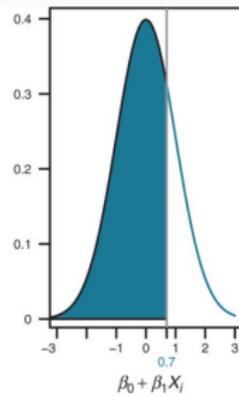
NORMAL DISTRIBUTION AND CDF

CDF = *Cumulative Density Function* of a normal distribution.

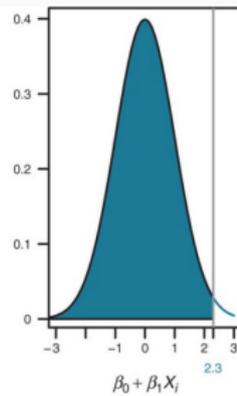
CDF written as probability of something being $<$ a certain value in a *normal distribution*. If we supply the CDF with any number it will give us the probability that a normally distributed random variable is less than that number.



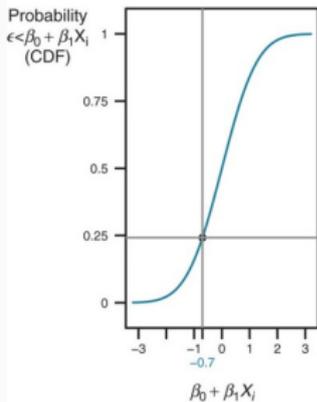
(a)



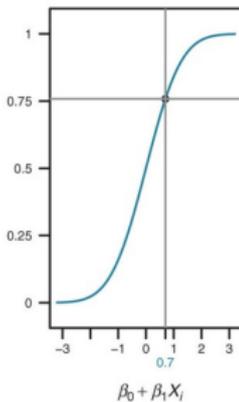
(b)



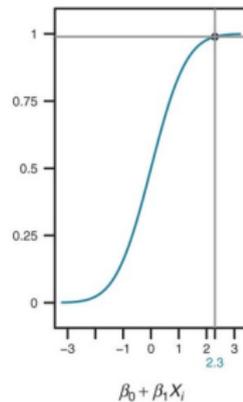
(c)



(d)



(e)



(f)

CDF has its own notation: the Greek letter Φ

CDF has its own notation: the Greek letter Φ indicates probability that a normally distributed random variable is less than whatever number we give it.

$$\begin{aligned} Pr(Y = 1) &= Pr(\epsilon \leq \beta_0 + \beta_1 X_1) \\ &= \Phi(\beta_0 + \beta_1 X_1) \end{aligned}$$

Probit estimates $\hat{\beta}$ that lead to high predicted probabilities (*close to 1*) for observations of Y that were actually **1**, and low predicted probabilities (*close to 0*) for observations of Y that were actually **0**.

Logit model is another way to estimate relationship between X and Y when the outcome is dichotomous.

Logit model is another way to estimate relationship between X and Y when the outcome is dichotomous.

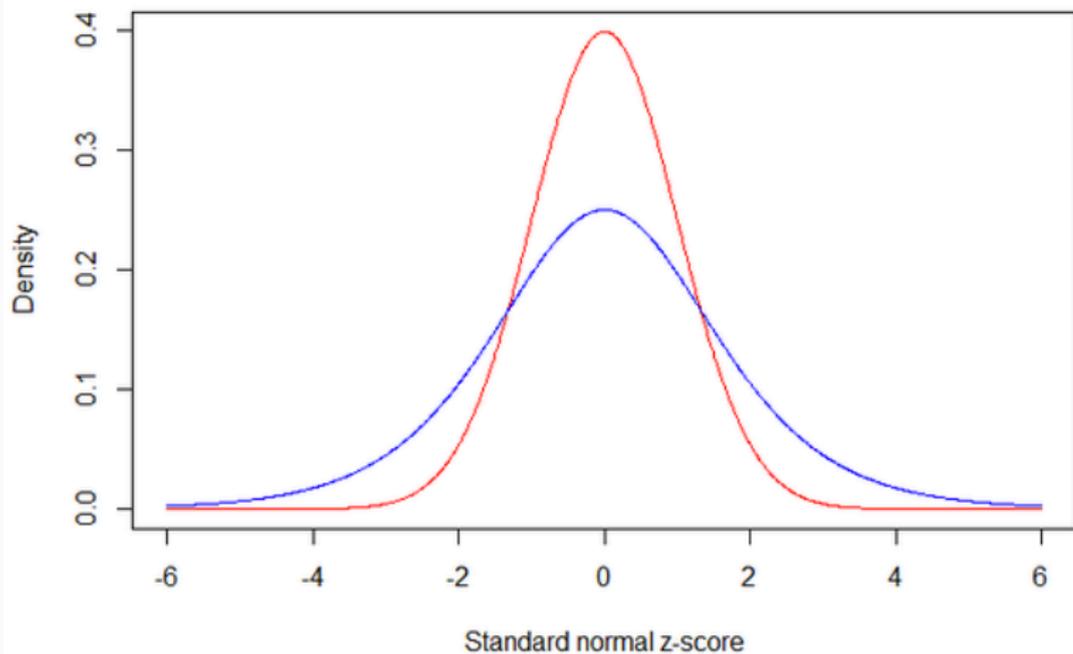
- The only mathematical difference between **logit** and **probit** is how we assume the error term is distributed...

Logit model is another way to estimate relationship between X and Y when the outcome is dichotomous.

- The only mathematical difference between **logit** and **probit** is how we assume the error term is distributed...
- In a logit model, we assume the error term follows a **logistic distribution**.

LOGISTIC VS. NORMAL DISTRIBUTION

Standard logistic vs. standard normal distribution



We again use **CDF** notation to specify how we estimate our model.

We again use **CDF** notation to specify how we estimate our model.

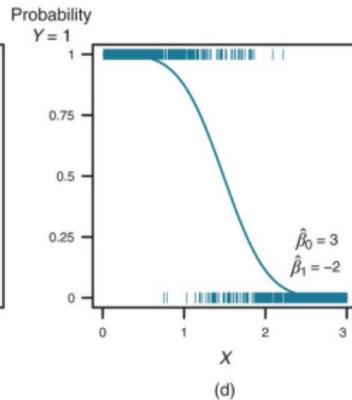
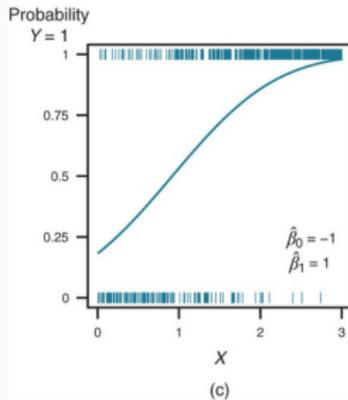
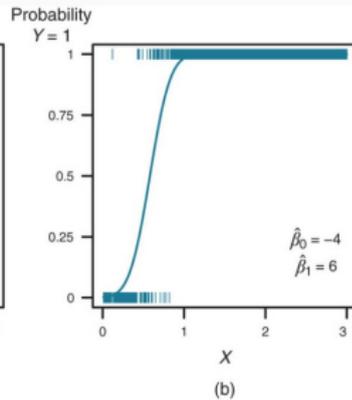
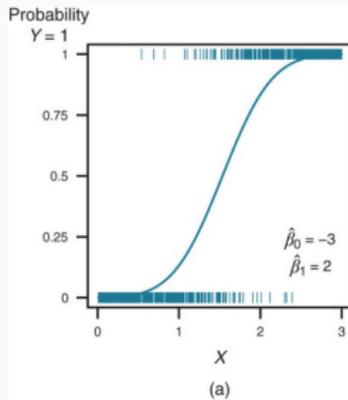
- When $\beta_0 + \beta_1 X_1$ is really large, logit function approaches **1**.

We again use **CDF** notation to specify how we estimate our model.

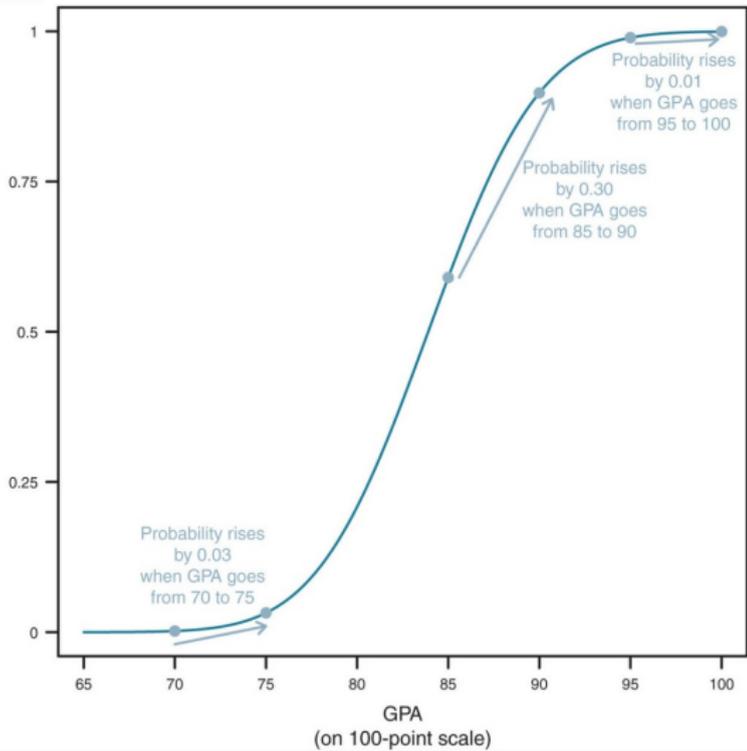
- When $\beta_0 + \beta_1 X_1$ is really large, logit function approaches **1**.
- When $\beta_0 + \beta_1 X_1$ is really small, logit function approaches **0**.

We again use **CDF** notation to specify how we estimate our model.

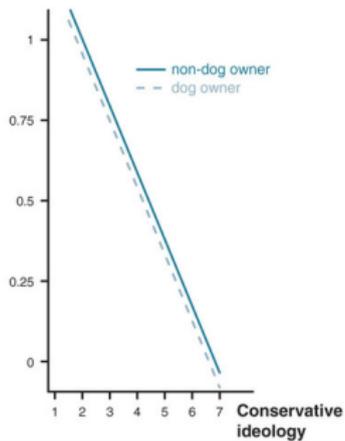
- When $\beta_0 + \beta_1 X_1$ is really large, logit function approaches **1**.
- When $\beta_0 + \beta_1 X_1$ is really small, logit function approaches **0**.
- When $\beta_0 + \beta_1 X_1$ is equal to 0, logit function is **0.5**.



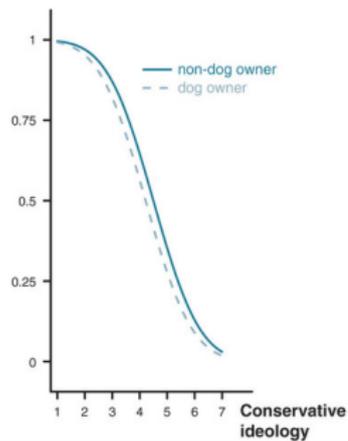
Probability of admission



Probability vote for Obama
Linear Probability Model (LPM)



Probability vote for Obama
Probit Model



Probability vote for Obama
Logit Model

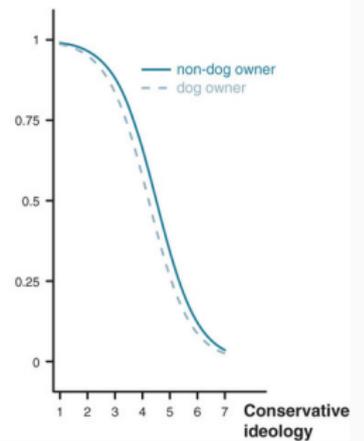


TABLE 12.3 Dog Ownership and Probability of Supporting Obama in 2008 Election

	LPM	Probit	Logit
Dog owner	-0.050* (0.006) [<i>t</i> = 8.14]	-0.213* (0.025) [<i>z</i> = 8.53]	-0.368* (0.044) [<i>z</i> = 8.44]
Ideology	-0.207* (0.002) [<i>t</i> = 106.0]	-0.753* (0.011) [<i>z</i> = 71.68]	-1.326* (0.021) [<i>z</i> = 63.98]
Constant	1.421* (0.009) [<i>t</i> = 159.9]	3.386* (0.049) [<i>z</i> = 69.74]	5.977* (0.095) [<i>z</i> = 62.84]
<i>N</i>	15,596	15,596	15,596
<i>R</i> ²	0.42		
log <i>L</i>		-6,720.87	-6,693.98
Minimum \hat{Y}_i	-0.080	0.018	0.024
Maximum \hat{Y}_i	1.213	0.995	0.990

Standard errors in parentheses.

** indicates significance at $p < 0.05$, two-tailed.*

TABLE 12.4 Estimated Effect of Dog Ownership and Ideology on Probability of Supporting Obama in 2008 Election

Variable	Simulated change	Probit	Logit
Dog owner	From 0 to 1, ideology at actual value	-0.052	-0.051
Ideology	Increase by 1, dog owner at actual value	-0.180	-0.181

DEALING WITH CATEGORICAL DV'S

Suppose we have data in *categories*, and we are sure of their *order*.

For example, we have this question from the ANES:

*There has been some discussion of abortion in recent years.
Which of the following best agrees with your view?*

ANES ABORTION QUESTION

1. By law, abortion should never be permitted.
2. The law should permit abortion only in case of rape, incest, or when the woman's life is in danger.
3. The law should permit abortion for reasons other than rape, incest, or danger to the woman's life, but only after the need for the abortion has been clearly established.
4. By law, a woman should always be able to obtain an abortion as a matter of personal choice.

HOW DO WE MODEL THIS?

- If we used OLS, we'd be assuming that each category is equally far apart.
- Instead, let's just assume that the range for each option is the same for everyone.
- This is the domain of **ordered logit/probit**.

- **Ordered** logit/probit take a **latent variable** approach.
- **Latent** means an unobserved tendency to have the observed outcome.
- We can't actually observe this underlying trait, but we can observe the outcome.

In ordered probit/logit, we assume that this unobserved y^* is:

- A function of X and whatever β we're estimating.
- Monotonically increasing (preserving order) in the value of the latent variable.

In ordered probit/logit, we assume that this unobserved y^* is:

- A function of X and whatever β we're estimating.
- Monotonically increasing (preserving order) in the value of the latent variable.

If the highest level of our outcome is 5 (“strongly agree”) on a 1-5 scale, then as y^* increases, our observed outcome should *never decrease*.

CATEGORIES AS BOUNDARIES

Values where the outcome “jumps” are **thresholds** (or boundaries or cutpoints): τ

CATEGORIES AS BOUNDARIES

Values where the outcome “jumps” are **thresholds** (or boundaries or cutpoints): τ

- Number of thresholds corresponds with the number of categories in our outcome variable (k).

$$\tau = k - 1$$

$$\tau = 5 - 1$$

$$\tau = 4$$

THRESHOLDS IN THE ABORTION EXAMPLE

For our abortion question, suppose we want to predict the answer with gender, ideology and income. We are thus saying:

$$y^* = \beta_1 \text{Female} + \beta_2 \text{Ideology} + \beta_3 \text{Income} + \epsilon$$

We don't observe y^* , so we model y .

We introduce three τ 's to separate our four categories.

INTERPRETING NON-LINEAR OUTPUT

INTERPRETATION

1. What is the quantity of interest?
2. What is the effect on Y of a small change in X ?
3. What is the prediction of Y given a profile of X ?
4. What is the difference in Y between two profiles of X ?

PARAMETERS IN NON-LINEAR MODELS

- In logit, the coefficient is a change in the “log-odds” of an outcome.
- Log-odds is not intuitive. You could convert to an odds-ratio. Still not intuitive.
- Can simply look at the sign of the coefficient, as a scan for substantive significance.
- Hard to compare coefficients because the effects of each depends on one others.

TWO MAIN WAYS

1. **Predicted probabilities:** probability of outcome being K given our X 's.
2. **Marginal effects:** discrete change associated with one unit increase in one covariate.

PREDICTED PROBABILITIES

- Predicted probabilities in R
- Predicted probs in logit and ordered logit
- Discrete marginal effects in R
- Marginal effects in R with cool plots
- Data camp practice course

REVIEWING PROBABILITY

AN INFORMATION PROBLEM

1. You have a question you want to ask about the world, or even a hunch about how things work, but you don't have complete information.

AN INFORMATION PROBLEM

1. You have a question you want to ask about the world, or even a hunch about how things work, but you don't have complete information.
2. We usually start somewhere, even with incomplete info (with a hunch, or a *prior*).

AN INFORMATION PROBLEM

1. You have a question you want to ask about the world, or even a hunch about how things work, but you don't have complete information.
2. We usually start somewhere, even with incomplete info (with a hunch, or a *prior*).
3. Our goal is to produce a **posterior probability estimate** (probability of observing evidence of your hypothesis given the data).

In the regression context, we're tempted to talk about β (or here, Θ) by asking:

What is the probability that the true β is or is not this big?

SHIFT FROM OLS

In the regression context, we're tempted to talk about β (or here, Θ) by asking:

What is the probability that the true β is or is not this big?

In reality, we have to say:

What is the probability, assuming that $\beta = 0$, that we'd observe an estimated $\hat{\beta}$ this far from 0?

Can't say anything about the true probability, but we can say something about relative probabilities.

1. Find the most likely Θ , given the data.
2. Compare different Θ 's, given the data.

FOR THURSDAY

- Class won't meet.
- Online tutorial and assignment instead.
- Homework will be up later today.