# Models of outcome and choice: The logit model

Silje Synnøve Lyder Hermansen

March 2, 2023

# Table of Contents

## Let's touch base

**We will be using mentimeter (menti.com) to communicate interactively.**

▶ answer questions on www.menti.com using the access code 8471 19241

▶ results show on screen

⇒ *Relax, your answers are anonymous!*

# Table of Contents

# Section 2

## A latent variable approach to GLMs

# Many outcomes are not continuous

**OLS assumes a continuous dependent variable. But many phenomena in the social sciences are not like that.**

▶ Vote choice, civil conflict onset, legislator performance, court rulings, time to compliance, etc.

▶ What phenomena are you interested in?

⇒ *OK. Let's strategize.*

# All regressions are linear(ized)

**The basic formulation in any regression describes a linear relationship between $x_i$ and $y_i$:**

$$y_i = \alpha + \beta x_i + \epsilon_i \tag{1}$$

- ▶ When $x_i$ increases with one unit, $y_i$ increases with $\beta$ units.
- ▶ If that relationship is not linear, we have to make it so:
    - ▶ by recoding the $x_i$
    - ▶ by recoding the $y_i$ → we *linearize*.

# A latent variable

**A linear(ized) model requires a continuious dependent variable.**

- ▶ Imagine we are interested in unobservable variable, $z_i$, that describes our propensity towards something.
- ▶ Above a certain threshold ($\tau$) of $z_i$, observability kicks in and we can see $y_i$.
- ▶ The regression coefficients ($\beta$) in GLMs describe that relationship.

$\Rightarrow$ The latent variable approach is useful when interpreting the results.

# Example: The binomial model

**The logit model is a perfect example:**

$$y_i = \begin{cases} 1 & \Leftrightarrow & z_i > \tau \\ 0 & \Leftrightarrow & z_i \leqslant \tau \end{cases} \tag{2}$$

▶ The probability $(z_i)$ of an outcome $y_i$ is continuous.

▶ Above a certain probability $(\tau)$, we observe a positive outcome $(y_i = 1)$.

$\Rightarrow$ *but how do we set the value of $\tau$?*

# From latent variable to descrete outcomes

**Statistical theory helps us describe how $z_i$ leads to $y_i$.**

▶ What kind of process generated our data? $\rightarrow$ data generating process (DGP)

▶ How can we best describe it? $\rightarrow$ choice of *probability distribution* (in GLM)

# The three components of GLMs

**When fitting the model, we need to make three choices:**

▶ A linear predictor: $\beta x_i$.

▶ A probability distribution: they're all in the exponential family

▶ A recoding strategy

# The three components of GLMs

**In R this translates to two additional arguments compared to your usual OLS.**

- ▶ A linear predictor: $\rightarrow$ (y $\sim$ x).
- ▶ A probability distribution: $\rightarrow$ (family =)
- ▶ A recoding strategy $\rightarrow$ (link = ).

# Latent variable approach for interpretation

▶ The latent variable approach is useful when interpreting results.
▶ That's when we map *from* the latent variable *to* the observed outcome.
⇒ *When estimating the model, we have to go the other way 'round.*

# Table of Contents

# Section 3

# Recoding: How do we get from a binary to a continuous variable?

## Data structure

**We can only observe the outcome produced by the latent variable. There are two data structures for binary data:**

▶ classes of observations: e.g.: rats in a cage, coin tosses...

▶ case-based: e.g.: legislator votes, Brexit...

## Data structure

**We can only observe the outcome produced by the latent variable. There are two data structures for binary data:**

- ▶ classes of observations: e.g.: rats in a cage, coin tosses... → *the closest to the latent continuous variable.*
- ▶ case-based: e.g.: legislator votes, Brexit...

⇒ *we know the number of successes and trials in a cage/class/stratum. That's our starting point.*

## Let's start with the odds

**Despite binary outcomes, we want a continuous variable that is unbounded at both ends. We define a stratum and start comparing**:

▶ Odds: Compare number of successes with number of failures within a stratum→ *continuous but highly skewed.*

▶ Logtransform the odds → *continuous and bell shaped.*

## Let's examplify with rats

**We kept a 1000 rats in a cage and a number of them died (failure) while others are still alive (success). How can we model this?**

# We calculate the odds

**We calculate the odds of surviving in a cage in a 1000 cages**

▶ Let's consider a series of 1000 trials where we let the successes go from complete failure (success = 0) to complete success (success = 1000)

```r
success = 0:1000
tries = 1000
#remember: failure = tries - success
odds <- success/(tries - success)

hist(odds, breaks = 100, col = "blue")

hist(log(odds), breaks = 101, col = "blue")

plot(log(odds), success, type = "l")
```
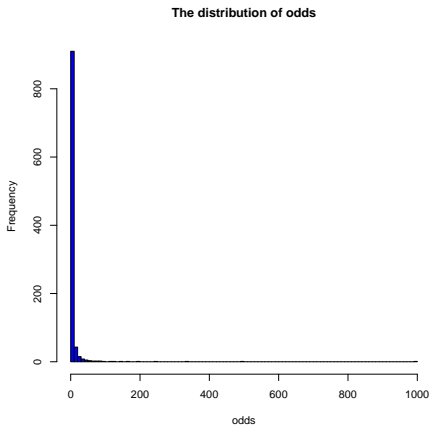
# Let's start with the odds

**We get a continuous but skewed variable.**



The distribution of odds

# Now, let's logtransform the odds

**We get a nice, bellshaped curve.**



The distribution of logodds

# Now, let's logtransform the odds

**This, we can run regressions on!**

## The famous S shape

**We can plot the logodds of success against the number of successes or their probability (it's the same).**



The relationship between probability and logodds

# Probability distributions for binary variables

**There are two, closely related probability distributions for binary outcomes**:
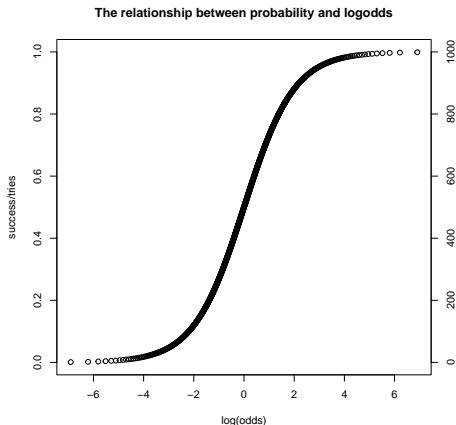
▶ The binomial distribution: $B(n, p)$

  ▶ $p$ is the probability of success tells where on the x-axis (trials) the distribution is placed.

  ▶ $n$ is the number of trials and defines the precision (width) of the distribution.

▶ The Bernoulli distribution: $Ber(p)$: when we only have only one trial.

Subsection 2

Why all the fuzz? Why not OLS?

# Distributions in OLS and maximum likelihood

▶ In OLS: The residuals must be normally distributed (but not the $y_i$)

▶ In ML: The $z_i$ must follow a known probability distribution.

⇒ *This what allows us to translate the latent variable to outcomes.*

# What happens if I run a linear model on binary outcomes?

▶ The model predicts out of the possible bounderies
  ▶ Predictions are wrong.
  ▶ Regression coefficients are wrong.
  ▶ Standard errors are wrong.

▶ The relationship between $x_i$ and $y_i$ is constant across all values.

⇒ *This last element has a bearing for the interpretation.*

# Table of Contents

# Section 4

## Interpretation: So... what did I find?

### Subsection 1

## Back and forth: Logistic and logit transformation

# The logit transformation

**When we go from outcomes to latent variable we use the logit transformation**.

$$logit(p) = log(\frac{p}{1-p}) \tag{3}$$

$\Rightarrow$ *This what R does when estimating our model*

# The logistic transformation

**When we go from the latent variable to outcomes we use the logistic transformation**.

$$logit^{-1}(logodds) = \frac{exp(logodds)}{1 + exp(logodds)} = \frac{1}{1 + exp(-logodds)} \quad (4)$$

$\Rightarrow$ *This what we do when interpreting our model*

# My three stages of interpretation

**I go through tree stages of interpretation**

- ▶ Inspect the marginal effects from regression table
    - ▶ Logodds: check direction and significance.
    - ▶ Odds ratio (for large coefficients) and percentage change (for smaller coefficients).
- ▶ Formulate scenarios using point estimates (in text)
- ▶ Formulate more scenarios with uncertainty using graphics.

```r
load("MEP2016.rda")
df <- MEP2016

mod <- glm(PoolsLocal ~
             OpenList +
             SeatsNatPal.prop +
             LaborCost,
           family = binomial(link = "logit"),
           df)

stargazer::stargazer(mod,
                     # label = "tab:regression",
                     title = "MEPs' propensity to share local
                     out = "results_table.tex",
                      type = "latex")
```

# The regression table: marginal effects

**I interpret the regression coefficient itself**

▶ Change in logodds: check direction and significance.

▶ Odds ratio (for large coefficients) and percentage change (for smaller coefficients).

⇒ *A first stab at hypothesis testing.*

# The regression table: marginal effects

**Now, you try!** What statements would you make using the change in logodds, the odds ratio and percentage change?

Table: MEPs' propensity to share local assistants (a binomial logit)

|  | Dependent variable: |
|---|---|
|  | PoolsLocal |
| OpenList | −1.124*** |
|  | (0.181) |
| SeatsNatPal.prop | −1.930*** |
|  | (0.527) |
| LaborCost | 0.056*** |
|  | (0.009) |
| Constant | −1.094*** |
|  | (0.286) |
| Observations | 686 |
| Log Likelihood | −392.832 |
| Akaike Inf. Crit. | 793.665 |
| *Note:* | *p<0.1; **p<0.05; ***p<0.01 |

```
#Change in logodds for MEPs in candidate-centered systems
mod$coefficients[2]

## OpenList
## -1.124427

# Odds ratio: <1 is negative; > 1 is positive
exp(mod$coefficients[2])

## OpenList
## 0.3248387

# Percentage change
(exp(mod$coefficients[2]) - 1)*100

## OpenList
## -67.51613
```

# The regression table: marginal effects

**Typical statements about marginal effects**

▶ Change in logodds: "MEPs from candidate-centered systems are less likely to share local assistants. Both effects are statistically significant."

▶ Percentage change (for smaller coefficients; -1.93)."The likelihood that an MEP shares a local assistant with a party colleague is 68% lower when they compete in a candidate-centered system compared to those that compete in party-centered systems."

⇒ *A first stab at hypothesis testing.*

# Predicted values

**If you believe the model describes reality appropriately, you can learn more about it by interpreting more thoroughly**

▶ Odds ratios are notoriously hard to understand.

▶ The effect depends on the value of $y_i$ and all the other $x$s.

$\Rightarrow$ *Interpret the predicted values*

# Predicted point estimates (text)

**Formulate scenarios using point estimates (in text)**

▶ Take an all-else-equal approach: Let one $x$ change and keep all others constant (on a typical value).

▶ Find the typical representative of two $x$ values and set the other $x$s accordingly.

$\Rightarrow$ *Which one you use depends on your objective: A theoretical point, assess effect of intervention on groups...*

# Predicted point estimates/first difference (text)

**Now you try!** What is the predicted effect of changing electoral system on MEPs' propensity to share local assistants ...

▶ In Bulgaria (Labor cost $==$ 4.4); when the party is small (Seat share $==$ 0.1).

▶ In Denmark (Labor cost $==$ 42); when the party is small (Seat share $==$ 0.1).

▶ Is this a realistic set of scenarios?

$\Rightarrow$ *Compare the two predicted probabilities for each pairs of scenarios.*

▶ Go to Padlet to provide your answer:
(https://padlet.com/siljesynnove/logit)

```r
dfp <- df %>% select(PoolsLocal,
                     SeatsNatPal.prop,
                     OpenList,
                     VoteShare_LastElection,
                     LaborCost)
```

```
## Error in df %>% select(PoolsLocal, SeatsNatPal.prop,
OpenList, VoteShare_LastElection, :  could not find
function "%>%"
```

```r
stargazer::stargazer(dfp,
                     title = "Descriptive statistics",
                     out = "desc_table.tex")
```

```
## Error in .stargazer.wrap(..., type = type, title =
title, style = style, :  object 'dfp' not found
```

```r
##Bulgaria; party-centered
logodds1 <- mod$coefficients[1] * 1 + mod$coefficients[2] * 0 +
```

```
  mod$coefficients[3] * 0.1 + mod$coefficients[4] * 4.4
prob1 <- exp(logodds1)/(1+exp(logodds1))
prob1

## (Intercept)
##   0.2610522

##Bulgaria; candidate-centered
logodds2 <- mod$coefficients[1] * 1 + mod$coefficients[2] * 1 +
  mod$coefficients[3] * 0.1 + mod$coefficients[4] * 4.4
prob2 <- exp(logodds2)/(1+exp(logodds2))
prob2

## (Intercept)
##    0.102944

#First difference
prob2 - prob1

## (Intercept)
##  -0.1581082

#Denmark; party-centered
logodds3 <- mod$coefficients[1] * 1 + mod$coefficients[2] * 0 +
  mod$coefficients[3] * 0.1 + mod$coefficients[4] * 42
prob3 <- exp(logodds3)/(1+exp(logodds3))
prob3

## (Intercept)
##   0.7430194

#Denmark; candidate-centered
logodds4 <- mod$coefficients[1] * 1 + mod$coefficients[2] * 1 +
  mod$coefficients[3] * 0.1 + mod$coefficients[4] * 42
prob4 <- exp(logodds4)/(1+exp(logodds4))
prob4
```

```
## (Intercept)
##   0.4843289

#First difference
prob4 - prob3

## (Intercept)
##  -0.2586905
```

```
#Alternative
newdata <- data.frame(OpenList = c(0, 1, 0, 1),
                      SeatsNatPal.prop = 0.1,
                      LaborCost = c(4.4, 4.4, 42, 42))

preds <- predict(mod, newdata, type = "response")
preds

##         1         2         3         4
## 0.2610522 0.1029440 0.7430194 0.4843289

diff(preds[1:2]); diff(preds[3:4])
```

```
##           2
## -0.1581082
##           4
## -0.2586905
```

# Predicted point estimates (text)

**Notice how the absolue effect of the electoral system changes!**

▶ **Marginal effect:** MEPs' likelihood of sharing assistants decreases by 68% % when we change electoral system. → *holds for all values of x.*

▶ **First difference (scenario 1a and b):** We see that changing electoral system when labor cost is *low* corresponds to a predicted 16 percentage points shift in likelihood of sharing assistants (from 26 percentage points to 10 percentage points).

▶ **First difference (scenario 2a and b):** We see that changing electoral system when labor cost is *high* corresponds to a predicted 26 percentage points shift in likelihood of sharing assistants (from 74 percentage points to 48 percentage points).

⇒ *This is an implicit interaction effect.*

# Predicted values (graphic)

**Formulate scenarios using point estimates and put them on speed**

▶ Predict $y$ values for the entire range of $x$ and plot it.

▶ Simulate confidence and plot that too.

▶ You can do this for two scenarios.

$\Rightarrow$ *You get a sense of the actual differences in the data.*

# Table of Contents

# Section 5

Model assessment: How well is reality described?

# Model assessment

**Model assessments aim to gauge how well we describe the data (i.e. the $y$).**

▶ comparison between predicted and observed values (as in OLS).

▶ mapping outcomes to the recoded, "latent" variable (GLM).

⇒ *You have a few additional "tricks" to the standard OLS assessment.*

# Brier score

**Describes the "average size" of the residuals.**

$$B_b \equiv \frac{1}{n}\Sigma_{i=1}^{n}(\hat{\theta}_i - y_i)^2 \tag{5}$$

$\Rightarrow$ *Lower scores imply better predictions.*

# How well do I discriminate?

**The real question for logits is how well do I distinguish 0s from 1s.**
⇒ *Several strategies.*

## Table comparison

**The real question for logits is how well do I distinguish 0s from 1s.**

▶ Table (e.g. $2 \times 2$) with proportion of predicted against observed values for 0s and 1s.

▶ It is $\chi^2$ distributed (ref. the Hosmer-Lemeshow test)

$\Rightarrow$ *But how do I set the cut values (the $\tau$)?*

# The ROC curve

**The ROC lets the cut values vary and displays how wrong we are on each side (true positive vs. false positive).**

▶ A model with good predictions has a curve tending towards the upper left corner.

▶ The actual cut value depends on our priorities

⇒ *The graphic is useful in and of itself*

# The separation plot

**The separation plot show how the density of observed "successes" increases as our predicted values increase.**
⇒ *Another graphic that is useful in and of itself*