

# Structural Models of Utility Maximization

Tyler Ransom

University of Oklahoma, Dept. of Economics

April 14, 2020

# Outline

## 1. Intro

## 2. Discrete choice

## 3. Math

## 4. Estimation

## 5. Sample selection

## 6. Dynamic Discrete Choice

# Today's plan

- 1 Describe static discrete choice models
- 2 How do they fit in with other data science models we've talked about in this class?
- 3 Derive logit/probit probabilities from intermediate microeconomic theory
- 4 Go through examples of how to estimate
- 5 How discrete choice models relate to sample selection bias

Note: These slides are based on the introductory lecture of a PhD course taught at Duke University by Peter Arcidiacono, and are used with permission. That course is based on Kenneth Train's book *Discrete Choice Methods with Simulation*, which is freely available [here](#) (PDF).

# Outline

1. Intro

**2. Discrete choice**

3. Math

4. Estimation

5. Sample selection

6. Dynamic Discrete Choice

# What are discrete choice models?

- ▶ Discrete choice models are one of the workhorses of structural economics
- ▶ Deeply tied to economic theory:
  - ▶ utility maximization
  - ▶ revealed preference
- ▶ Used to model “utility” (broadly defined), for example:
  - ▶ consumer product purchase decisions
  - ▶ firm market entry decisions
  - ▶ investment decisions

## Why use discrete choice models?

- ▶ Provides link between human optimization behavior and economic theory
- ▶ Parameters of these models map directly to economic theory
- ▶ Parameter values can quantify a particular policy
- ▶ Can be used to form counterfactual predictions (e.g. by adjusting certain parameter values)
- ▶ Allows a research to quantify “tastes”

# Why not use discrete choice models?

- ▶ They're not the best predictive models
  - ▶ Trade-off between out-of-sample prediction and counterfactual prediction
- ▶ You don't want to form counterfactual predictions, you just want to be able to predict handwritten digits
- ▶ You aren't interested in economic theory
- ▶ The math really scares you
- ▶ You don't like making assumptions
  - ▶ e.g. that decision-makers are rational

## Example of a discrete choice model

- ▶ Cities in the Bay Area are interested in how the introduction of rideshare services will impact ridership on Bay Area Rapid Transit (BART)
- ▶ Questions that cities need to know the answers to:
  - ▶ Is rideshare a substitute for public transit or a complement?
  - ▶ How inelastic is demand for BART? Should fares be  $\uparrow$  or  $\downarrow$ ?
  - ▶ Should BART services be scaled up to compete with rideshares?
  - ▶ Will the influx of rideshare vehicles increase traffic congestion / pollution?
- ▶ Each of these questions requires making a counterfactual prediction
- ▶ In particular, need a way to make such a prediction confidently and in a way that is easy to understand



# Properties of discrete choice models

- 1 Agents choose from among a **finite** set of alternatives (called the *choice set*)
- 2 Alternatives in choice set are **mutually exclusive**
- 3 Choice set is **exhaustive**

## Example illustrating these properties

- ▶ In San Francisco, people can commute to work by the following (and *only* the following) methods:
  - ▶ Drive a personal vehicle (incl. motorcycle)
  - ▶ Carpool in a personal vehicle
  - ▶ Use taxi/rideshare service (incl. Uber, Lyft, UberPool, LyftLine, etc.)
  - ▶ BART (bus, train, or both)
  - ▶ Bicycle
  - ▶ Walk

# Outline

1. Intro

2. Discrete choice

**3. Math**

4. Estimation

5. Sample selection

6. Dynamic Discrete Choice

## Mathematically representing utility

Let  $d_i$  indicate the choice individual (or decision-maker)  $i$  makes where  $d_i \in \{1, \dots, J\}$ . Individuals choose  $d$  to maximize their utility,  $U$ .  $U$  generally is written as:

$$U_{ij} = u_{ij} + \varepsilon_{ij} \quad (1)$$

where:

- 1  $u_{ij}$  relates observed factors to the utility individual  $i$  receives from choosing option  $j$
- 2  $\varepsilon_{ij}$  are unobserved to the researcher but observed to the individual
- 3  $d_{ij} = 1$  if  $u_{ij} + \varepsilon_{ij} > u_{ij'} + \varepsilon_{ij'}$  for all  $j' \neq j$

## Breakdown of the assumptions

- ▶ Examples of what's in  $\varepsilon$ 
  - ▶ Person's mental state when making the decision
  - ▶ Choices of friends or relatives (maybe, depends on the data)
  - ▶  $\vdots$
  - ▶ Anything else about the person that is not in our data
- ▶ Reasonable to assume additive separability?
  - ▶ This is a big assumption: that there are no interactive effects between unobservable and observable factors
  - ▶ This results in linear separation regions and may be too restrictive
  - ▶ For now, go with it, and remember that there are no free lunches

## Probabilistic choice

With the  $\varepsilon$ 's unobserved, we must consider choices as probabilistic instead of certain. The Probability that  $i$  chooses alternative  $j$  is:

$$P_{ij} = \Pr(u_{ij} + \varepsilon_{ij} > u_{ij'} + \varepsilon_{ij'} \quad \forall j' \neq j) \quad (2)$$

$$= \Pr(\varepsilon_{ij'} - \varepsilon_{ij} < u_{ij} - u_{ij'} \quad \forall j' \neq j) \quad (3)$$

$$= \int_{\varepsilon} I(\varepsilon_{ij'} - \varepsilon_{ij} < u_{ij} - u_{ij'} \quad \forall j' \neq j) f(\varepsilon) d\varepsilon \quad (4)$$

## Transformations of utility

Note that, regardless of what distributional assumptions are made on the  $\varepsilon$ 's, the probability of choosing a particular option does not change when we:

- ➊ Add a constant to the utility of all options (utility is relative to one of the options, only differences in utility matter)
- ➋ Multiply by a positive number (need to scale something, generally the variance of the  $\varepsilon$ 's)

This is just like in consumer choice theory: utility is ordinal, and so is invariant to the above two transformations

## Variables

Suppose we have:

$$u_{i1} = \alpha \text{Male}_i + \beta_1 X_i + \gamma Z_1$$

$$u_{i2} = \alpha \text{Male}_i + \beta_2 X_i + \gamma Z_2$$

Since only differences in utility matter:

$$u_{i1} - u_{i2} = (\beta_1 - \beta_2)X_i + \gamma(Z_1 - Z_2)$$

- ▶ Thus, we cannot tell whether men are happier than women, but can tell whether men have a preference for a particular option over another.
- ▶ We can only obtain **differenced** coefficient estimates on  $X$ 's, and can obtain an estimate of a coefficient that is constant across choices only if the variable it is multiplying varies by choice.



## Number of error terms

Similar to socio-demographic characteristics, there are restrictions on the number of error terms. Recall that the probability  $i$  will choose  $j$  is given by:

$$\begin{aligned}P_{ij} &= \Pr(u_{ij} + \varepsilon_{ij} > u_{ij'} + \varepsilon_{ij'} \quad \forall j' \neq j) \\&= \Pr(\varepsilon_{ij'} - \varepsilon_{ij} < u_{ij} - u_{ij'} \quad \forall j' \neq j) \\&= \int_{\varepsilon} I(\varepsilon_{ij'} - \varepsilon_{ij} < u_{ij} - u_{ij'} \quad \forall j' \neq j) f(\varepsilon) d\varepsilon\end{aligned}$$

where the integral is  $J$ -dimensional.

## Number of error terms (cont'd)

But we can rewrite the last line as  $J - 1$  dimensional integral over the differenced  $\varepsilon$ 's:

$$P_{ij} = \int_{\tilde{\varepsilon}} I(\tilde{\varepsilon}_{ij'} < \tilde{u}_{ij'} \quad \forall j' \neq j) g(\tilde{\varepsilon}) d\tilde{\varepsilon}$$

Note that this means one dimension of  $f(\varepsilon)$  is not identified and must therefore be normalized.

## Derivation of Logit Probability

Consider the case when the choice set is  $\{1, 2\}$ . The Type 1 extreme value cdf for  $\varepsilon_2$  is:

$$F(\varepsilon_2) = e^{-e^{(-\varepsilon_2)}}$$

To get the probability of choosing 1, substitute in for  $\varepsilon_2$  with  $\varepsilon_1 + u_1 - u_2$ :

$$Pr(d_1 = 1 | \varepsilon_1) = e^{-e^{-(\varepsilon_1 + u_1 - u_2)}}$$

But  $\varepsilon_1$  is unobserved so we need to integrate it out (see Appendix to these slides if you want the math steps)

## Derivation of Logit Probability

In the end, we can show that, for any model where there are two choice alternatives and  $\varepsilon$  is drawn from the Type 1 extreme value distribution,

$$P_{i1} = \frac{\exp(u_{i1} - u_{i2})}{1 + \exp(u_{i1} - u_{i2})}, P_{i2} = \frac{1}{1 + \exp(u_{i1} - u_{i2})}$$

Suppose we have a data set with  $N$  observations. The log likelihood function we maximize is then:

$$\ell(\beta, \gamma) = \sum_{i=1}^N (d_{i1} = 1)(u_{i1} - u_{i2}) - \ln(1 + \exp(u_{i1} - u_{i2}))$$

## Derivation of Probit Probability

In the probit model, we assume that  $\varepsilon$  is Normally distributed. So for a binary choice we have:

$$P_{i1} = \Phi(u_{i1} - u_{i2}), P_{i2} = 1 - \Phi(u_{i1} - u_{i2})$$

where  $\Phi(\cdot)$  is the standard normal cdf

The log likelihood function we maximize is then:

$$\ell(\beta, \gamma) = \sum_{i=1}^N (d_{i1} = 1) \ln(\Phi(u_{i1} - u_{i2})) + (d_{i2} = 1) \ln(1 - \Phi(u_{i1} - u_{i2}))$$

## Pros & Cons of Logit & Probit

Logit model:

- ▶ Has a much simpler objective function
- ▶ Is by far most popular
- ▶ ... but has more restrictive assumptions about how people substitute choices
- ▶ (this is known as the Independence of Irrelevant Alternatives or IIA assumption)

Probit model:

- ▶ Much more difficult to estimate
- ▶ ... but can accommodate more realistic choice patterns

# Outline

1. Intro

2. Discrete choice

3. Math

**4. Estimation**

5. Sample selection

6. Dynamic Discrete Choice

## Estimation in R

The R function `glm` is the easiest way to estimate a binomial logit or probit model:

```
library(mlogit)
data(Heating) # load data on residential heating choice in CA
levels(Heating$depvar) <- c("gas","gas","elec","elec","elec")
estim <- glm(depvar ~ income+agehed+rooms+region,
               family=binomial(link='logit'),data=Heating))
print(summary(estim))
```



# Interpreting the coefficients

Estimated coefficients using the code in the previous slide:

Coefficients:

	Estimate	Std. Error	z value	Pr(> z )
(Intercept)	-0.599142	0.252032	-2.377	0.0174 *
income	0.015579	0.027905	0.558	0.5766
agehed	-0.006535	0.003342	-1.955	0.0505 .
rooms	0.024291	0.026916	0.902	0.3668
regionscostl	-0.053096	0.126665	-0.419	0.6751
regionmountn	0.041827	0.169787	0.246	0.8054
regionncostl	-0.219136	0.137692	-1.591	0.1115

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 '\_' 1

## Interpreting the coefficients

- ▶ Positive coefficients  $\Rightarrow$  household more likely to choose the non-baseline alternative (in this case: electric)
  - ▶ Whatever the first level of the factor dependent variable is will be the "baseline" alternative
- ▶ Negative coefficients imply the reverse
- ▶ Coefficients **not** linked to changes in probability of choosing the alternative (since probability is a nonlinear function of  $X$ )

## Forming predictions

To get predicted probabilities for each observation in the data:

```
Heating$predLogit <- predict(estim, newdata = Heating, type = "response")  
print(summary(Heating$predLogit))
```

## Estimating a probit model

For the probit model, we repeat the same code, except change the “link” function from “logit” to “probit”

```
estim2 <- glm(depvar ~ income+agehed+rooms+region,  
              family=binomial(link='probit'), data=Heating))  
print(summary(estim2))  
Heating$predProbit <- predict(estim2, newdata = Heating, type = "response")  
print(summary(Heating$predProbit))
```

## A simple counterfactual simulation

- ▶ We talked a lot about doing counterfactual comparisons, but how do we *actually* do it?
- ▶ Let's show how to do this on a previous example. Suppose that we introduce a policy that makes richer people more likely to use electric heating.
- ▶ Mathematically, what does this look like?
- ▶ It would correspond to an increase in the parameter in front of *income* in our regression

## A simple counterfactual simulation

- Suppose the coefficient increased by a factor of 4. What is the new share of gas vs. electricity usage?

```
estim$coefficients["income"] <- 4*estim$coefficients["income"]  
Heating$predLogitCfl <- predict(estim, newdata = Heating, type = "  
    response")  
print(summary(Heating$predLogitCfl))
```

This policy would increase electric usage by 7 percentage points (from 22% to 29%)

# Outline

1. Intro
2. Discrete choice
3. Math
4. Estimation
- 5. Sample selection**
6. Dynamic Discrete Choice

# Discrete choice models and sample selection bias

- ▶ Discrete choice models are common tools used to evaluate sample selection bias
- ▶ Why? Because variables that are MNAR can be thought of as following a utility-maximizing process
- ▶ Examples:
  - ▶ Suppose you want to know what the returns to schooling are, but you only observe wages for those who currently hold jobs
  - ▶ As a result, your estimate of the returns to schooling might be invalidated by the non-randomness of the sample of people who are currently working
  - ▶ How to get around this? Use a discrete choice model (This was the problem we ran into in PS7, if you recall)



## Heckman selection correction

The Heckman selection model specifies two equations:

$$u_i = \beta x_i + \nu_i$$

$$y_i = \gamma z_i + \varepsilon_i$$

- ▶ The first equation is a utility maximization problem, determining if the person is in the labor force. Can think of  $\nu_i$  as “desire to work”
- ▶  $x_i$  may include: number of children in the household
- ▶ The second equation is the log wage equation, where  $y_i$  is only observed for people who are in the labor force.
- ▶ To solve the model, one needs to use the so-called “Heckit” model, which involves adding a correction term in the wage equation which accounts for the fact that workers are not randomly selected.

## Estimating Heckman selection in R

R has a package called `sampleSelection` which incorporates the Heckman selection model<sup>1</sup>

```
library(sampleSelection)
data('Mroz87')
Mroz87$kids <- (Mroz87$kids5 + Mroz87$kids618) > 0
# Comparison of linear regression and selection model
outcome1 <- lm(wage ~ exper, data = Mroz87)
summary(outcome1)
selection1 <- selection(selection = lfp ~ age + I(age^2) + faminc + kids
  + educ,
outcome = wage ~ exper, data = Mroz87, method = '2step')
summary(selection1)
```

---

<sup>1</sup>This code taken from Garrett Glasgow's website:

[http://www.polsci.ucsb.edu/faculty/glasgow/ps207/ps207\\_class6.r](http://www.polsci.ucsb.edu/faculty/glasgow/ps207/ps207_class6.r)

## Estimation output

Output from a regression of wage on experience:

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	1.30434	0.18937	6.888	1.20e-11	***
exper	0.10067	0.01419	7.093	3.03e-12	***
- - -					

## Estimation output

Output from the Heckman selection model: (edited for length)

Probit selection equation:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-4.157e+00	1.402e+00	-2.965	0.003126	**
kidsTRUE	-4.490e-01	1.309e-01	-3.430	0.000638	***

Outcome equation:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	7.12492	0.80425	8.859	<2e-16	***
exper	0.02962	0.02059	1.439	0.151	

Multiple R-Squared:0.0823, Adjusted R-Squared:0.0779

Error terms:

	Estimate	Std. Error	t value	Pr(> t )	
invMillsRatio	-5.075	1.108	-4.581	5.42e-06	***
sigma	4.977	NA	NA	NA	
rho	-1.020	NA	NA	NA	

## Reading the output

- ▶ Because there are two equations, there are now more parameters
- ▶ Using just the regression on workers led us to believe the returns to experience were  $\approx 10\%$
- ▶ Taking into account the selectivity of labor force participants leads us to conclude the returns to experience are much lower ( $\approx 3\%$  and not statistically different from zero)
- ▶ Viability of the model depends on the assumption that's made: in this case, that having children only affects labor supply preferences and doesn't affect wages
  - ▶ Wage discrimination against mothers in the labor market would invalidate this assumption
  - ▶ Back to the idea that to get causal inference we have to impose more assumptions

# Outline

1. Intro

2. Discrete choice

3. Math

4. Estimation

5. Sample selection

**6. Dynamic Discrete Choice**

# The optimal stopping problem

- ▶ Much of life is concerned with knowing when to stop:
  - ▶ How many people to date before making/accepting a marriage proposal
  - ▶ How much to study for upcoming exams
  - ▶ How long to “hodl” an asset
- ▶ All of the above cases involve forming expectations about:
  - 1 The long-run value of making a particular choice
  - 2 ... relative to the long-run value of alternatives
- ▶ Expectations about the future imply that we need to think “dynamically” (i.e. think over the longterm)
- ▶ Today we'll go through the math on how to do this

## Relation to reinforcement learning

- ▶ Reinforcement learning is based on the optimal stopping problem
- ▶ At each state  $X$  (e.g. game board configuration), observe reward  $y$  (e.g. win probability)
- ▶ In each period (i.e. gameplay turn), choose the decision that maximizes the (present value) expected reward
- ▶ With structural models, “reward” is utility



## Dynamic discrete choice models

With *dynamic* models, need a way to quantify present value of utility  
Individual  $i$ 's **flow utility** for option  $j$  at time  $t$  is:

$$\begin{aligned} U_{ijt} &= u_{ijt} + \varepsilon_{ijt} \\ &= X_{it}\alpha_j + \varepsilon_{ijt} \end{aligned}$$

Individual chooses  $d_{it}$  to maximize **expected lifetime utility**

$$\max_{d_{it}} V = E \left\{ \sum_{\tau=t}^T \sum_j \beta^{\tau-t} (d_{it} = j) U_{ijt} \right\}$$

- ▶  $V$  is the *value function*
- ▶  $\beta \in (0, 1)$  is the *discount factor*
- ▶  $T$  is the *time horizon*

# Expectations

- ▶ Expectations taken over future states ( $X$ 's) **and** errors
- ▶  $\varepsilon$ 's are iid over time
- ▶ Future states are not affected by  $\varepsilon$ 's except through current and past choices:

$$E(X_{t+1} | d_t, \dots, d_1, \varepsilon_t, \dots, \varepsilon_1) = E(X_{t+1} | d_t, \dots, d_1)$$

# Human behavior vs. reinforcement learning

- ▶ In reinforcement learning, we typically don't have  $\varepsilon$ , unless we want to allow for "curiosity"
- ▶ Transitions in  $X$  much more dominant factor (e.g. if I move here, opponent will move there, ...)
- ▶ Real-life example of uncertainty in  $\varepsilon$ 's:
  - ▶ "My significant other might take a job in another city next year, so if I want to move with him/her, I may not want to take this job offer today."
- ▶ Real-life example of uncertainty in  $X$ 's:
  - ▶ "I might get laid off next year, which will influence my ability to pay off my car loan, so I might want not want to buy this Mercedes today, since my (expected) permanent income might be lower than my current income."

# Dynamic programming & the Bellman equation

- ▶ We want to maximize the value function  $V$
- ▶ It's helpful to write the value function as a recursive expression, where we separate out today's decision from all future decisions (this is called the *Bellman equation*, or the *dynamic programming problem*)
- ▶ The payoff from choosing alternative  $j$  today is the *flow utility*  $= u_{ijt}$  from earlier in these slides
- ▶ The payoff from choosing alternative  $j$  in the future is the expected future utility conditional on choosing  $j$  today

How do we solve the Bellman equation?

- ▶ Requires solving backwards, just like in a dynamic game (cf. subgame perfect Nash equilibrium)

## Two Period Example

Consider the utility of choice  $j$  in the last period:

$$\begin{aligned} U_{ijT} &= u_{ijT} + \varepsilon_{ijT} \\ &= X_{iT}\alpha_j + \varepsilon_{ijT} \end{aligned}$$

Define the **conditional valuation function** for choice  $j$  as the flow utility of  $j$  minus the associated  $\varepsilon$  plus the expected value of future utility conditional on  $j$ :

$$v_{ijT-1} = u_{ijT-1} + \beta E \max_{k \in J} \{u_{ikT} + \varepsilon_{ikT} | d_{iT-1} = j\}$$

where  $\beta$  is the discount factor.

Suppose  $X_{iT}$  was deterministic given  $X_{iT-1}$  and  $d_{T-1}$  and the  $\varepsilon$ 's are Type 1 extreme value. What would the  $E \max$  expression be?  $[\ln \sum_k \exp(u_{ikT})]$

## Two Period Example (cont'd)

For  $J = 2$  the log likelihood would then look like:

$$L(\alpha) = \sum_{i=1}^N \sum_{t=1}^T (d_{i1t} = 1)(v_{i1t} - v_{i2t}) - \ln(1 + \exp(v_{i1t} - v_{i2t}))$$

where

$$v_{ijt} = u_{ijt} + \beta E \max_{k \in J} \{v_{ikt+1} + \varepsilon_{ikt+1} | d_{it} = j\}$$

and where

$$u_{ijt} = X_{it}\alpha_j$$

Note: if  $T = 2$  then  $v_{ikt+1} = u_{ikT}$

## Estimating a dynamic discrete choice model in R

- ▶ Because we have to loop backwards through time, we can't use a canned function like `lm()`
- ▶ Requires us to write a custom likelihood function
- ▶ This is because the flow utility parameters ( $\alpha_j$ ) appear in the flow utility function in *each* period
  - ▶ Side note: We don't typically estimate the discount factor ( $\beta$ ) but instead assume a fixed value (most common: 0.90 or 0.95)
- ▶ To do this, write down an objective function (i.e. log likelihood function) and use `nloptr` to estimate the  $\alpha$ 's
- ▶ Once you have the  $\alpha$ 's you can do counterfactual simulations
- ▶ These simulations are likely to be more realistic because the model has incorporated forward-looking behavior

## Objective function

```
objfun <- function(alpha,Choice,age) {  
  J <- 2  
  a <- alpha[3]*(1-diag(J))  
  
  u1 <- matrix(0, N, T)  
  u2 <- matrix(0, N, T)  
  for (t in 1:T) {  
    u1[ ,t] <- 0*age[ ,t]  
    u2[ ,t] <- alpha[1] + alpha[2]*age[ ,t]  
  }  
}
```

(continued on next slide)



```

Like <- 0
for (t in T:1) {
  for (j in 1:J) {
    # Generate FV
    dem <- exp(u1[ ,t] + a[1,j]+fv[ ,1,t+1])+
           exp(u2[ ,t] + a[2,j]+fv[ ,2,t+1])
    fv[ ,j,t] <- beta*(log(dem)-digamma(1))
    p1 <- exp(u1[ ,t] + a[1,j] + fv[ ,1,t+1])/dem
    p2 <- exp(u2[ ,t] + a[2,j] + fv[ ,2,t+1])/dem
    Like <- Like - (LY[ ,t]==j)*((Choice[ ,t]==1)*log(p1)+(Choice
      [ ,2]==2)*log(p2))
  }
}
return ( sum(Like) )
}

```

## Calling nloptr

```
## initial values
theta0 <- runif(3) #start at uniform random numbers equal to number of
  coefficients

## Algorithm parameters
options <- list("algorithm"="NLOPT_LN_NELDERMEAD","xtol_rel"=1.0e-6,"
  maxeval"=1e4)

## Optimize!
result <- nloptr( x0=theta0,eval_f=objfun,opts=options,Choice=Choice,age=
  age)
print(result)
```

# Derivation of Logit Probability

$$\begin{aligned}Pr(d_1 = 1) &= \int_{-\infty}^{\infty} \left( e^{-e^{-(\varepsilon_1 + u_1 - u_2)}} \right) f(\varepsilon_1) d\varepsilon_1 \\&= \int_{-\infty}^{\infty} \left( e^{-e^{-(\varepsilon_1 + u_1 - u_2)}} \right) e^{-\varepsilon_1} e^{-e^{-\varepsilon_1}} d\varepsilon_1 \\&= \int_{-\infty}^{\infty} \exp \left( -e^{-\varepsilon_1} - e^{-(\varepsilon_1 + u_1 - u_2)} \right) e^{-\varepsilon_1} d\varepsilon_1 \\&= \int_{-\infty}^{\infty} \exp \left( -e^{-\varepsilon_1} [1 + e^{u_2 - u_1}] \right) e^{-\varepsilon_1} d\varepsilon_1\end{aligned}$$

# Derivation of Logit Probability

Now need to do the substitution rule where  $t = \exp(-\varepsilon_1)$  and  $dt = -\exp(-\varepsilon_1)d\varepsilon_1$ .

Note that we need to do the same transformation of the bounds as we do to  $\varepsilon_1$  to get  $t$ . Namely,  $\exp(-\infty) = 0$  and  $\exp(\infty) = \infty$ .

# Derivation of Logit Probability

Substituting in then yields:

$$\begin{aligned} Pr(d_1 = 1) &= \int_{-\infty}^0 \exp(-t[1 + e^{(u_2 - u_1)}]) (-dt) \\ &= \int_0^{\infty} \exp(-t[1 + e^{(u_2 - u_1)}]) dt \\ &= \frac{\exp(-t[1 + e^{(u_2 - u_1)}])}{-[1 + e^{(u_2 - u_1)}]} \bigg|_0^{\infty} \\ &= 0 - \frac{1}{-[1 + e^{(u_2 - u_1)}]} = \frac{\exp(u_1)}{\exp(u_1) + \exp(u_2)} \end{aligned}$$