

# On the initialization of multinomial probit models

Lennart Oelschläger

Bielefeld University  
Department of Business Administration and Economics  
Econometrics Group of Prof. Dr. Bauer

July 16, 2021

## 1 The multinomial probit model

## 2 Initialization



### Description

This function carries out a minimization of the function  $f$  using a Newton-type algorithm. See the references for details.

### Usage

```
nlm(f, p, ..., hessian = FALSE, tysize = rep(1, length(p)),  
     fscale = 1, print.level = 0, ndigit = 12, gradtol = 1e-6,  
     stepmax = max(1000 * sqrt(sum((p/tysize)^2)), 1000),  
     steptol = 1e-6, iterlim = 100, check.analyticals = TRUE)
```

### Arguments

$f$

the function to be minimized, returning a single numeric value. This should be a function with first arguments specified by the ... argument.

If the function value has an attribute called `gradient` or both `gradient` and `hessian` attributes, these will be used. Otherwise, numerical derivatives are used. [deriv](#) returns a function with suitable `gradient` attribute and

$p$

starting parameter values for the minimization.



## 1 The multinomial probit model

- Definition
- Parameters
- Likelihood

## 2 Initialization

## The multinomial probit model

$$\overbrace{\begin{pmatrix} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{pmatrix}}^X$$

## The multinomial probit model

$$\begin{array}{c} \overbrace{\left( \begin{array}{ccc} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{array} \right)}^X \quad \overbrace{\left( \begin{array}{c} \beta_1 \\ \vdots \\ \beta_P \end{array} \right)}^\beta \end{array}$$

## The multinomial probit model

$$\overbrace{\begin{pmatrix} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{pmatrix}}^X \quad \overbrace{\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_P \end{pmatrix}}^\beta \quad + \quad \overbrace{\begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_J \end{pmatrix}}^\epsilon$$

## The multinomial probit model

$$\begin{matrix} \underbrace{U} \\ \begin{pmatrix} U_1 \\ \vdots \\ U_J \end{pmatrix} \end{matrix} = \begin{matrix} \underbrace{X} \\ \begin{pmatrix} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{pmatrix} \end{matrix} \begin{matrix} \underbrace{\beta} \\ \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_P \end{pmatrix} \end{matrix} + \begin{matrix} \underbrace{\epsilon} \\ \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_J \end{pmatrix} \end{matrix}$$



## The multinomial probit model

$$\begin{matrix} \underbrace{U} \\ \begin{pmatrix} U_1 \\ \vdots \\ U_J \end{pmatrix} \end{matrix} = \begin{matrix} \underbrace{X} \\ \begin{pmatrix} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{pmatrix} \end{matrix} \begin{matrix} \underbrace{\beta} \\ \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_P \end{pmatrix} \end{matrix} + \begin{matrix} \underbrace{\epsilon} \\ \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_J \end{pmatrix} \end{matrix}$$

$$y = \arg \max U$$

## The multinomial probit model

$$\underbrace{\begin{pmatrix} U_1 \\ \vdots \\ U_J \end{pmatrix}}_U = \underbrace{\begin{pmatrix} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta_1 \\ \vdots \\ \beta_P \end{pmatrix}}_\beta + \underbrace{\begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_J \end{pmatrix}}_\epsilon$$

$$y = \arg \max U$$

$$\epsilon = L \underbrace{\eta}_{\sim N_J(0, I)} \sim N_J(0, LL' = \Sigma)$$

## The multinomial probit model

$$\underbrace{\begin{pmatrix} U \\ U_1 \\ \vdots \\ U_J \end{pmatrix}}_U = \underbrace{\begin{pmatrix} X_{11} & \dots & X_{1P} \\ \vdots & \ddots & \vdots \\ X_{J1} & \dots & X_{JP} \end{pmatrix}}_X \underbrace{\begin{pmatrix} \beta \\ \beta_1 \\ \vdots \\ \beta_P \end{pmatrix}}_\beta + \underbrace{\begin{pmatrix} \epsilon \\ \epsilon_1 \\ \vdots \\ \epsilon_J \end{pmatrix}}_\epsilon$$

$$y = \arg \max U$$

$$\epsilon = L \underbrace{\eta}_{\sim N_J(0, I)} \sim N_J(0, LL' = \Sigma)$$

$$\beta = b + O \underbrace{\eta}_{\sim N_P(0, I)} \sim N_P(b, OO' = \Omega)$$

## Level normalization

$$\Delta_J U = \Delta_J X \beta + \Delta_J \epsilon, \quad \Delta_J = (I_{J-1} \quad -1) \in \mathbb{R}^{(J-1) \times J}$$

$$y = \begin{cases} i, & U_i = \max \Delta_J U > 0, i = 1, \dots, J-1 \\ J, & \Delta_J U < 0 \end{cases}$$

$$\Delta_J \epsilon \sim N_{J-1}(0, \Delta_J L L' \Delta_J')$$

## Level normalization

$$\Delta_J U = \Delta_J X \beta + \Delta_J \epsilon, \quad \Delta_J = (I_{J-1} \quad -1) \in \mathbb{R}^{(J-1) \times J}$$

$$y = \begin{cases} i, & U_i = \max \Delta_J U > 0, i = 1, \dots, J-1 \\ J, & \Delta_J U < 0 \end{cases}$$

$$\Delta_J \epsilon \sim N_{J-1}(0, \Delta_J L L' \Delta_J')$$

## Scale normalization

$$(\Delta_J L)_{11} = 1$$

## Number of parameters to estimate

alternatives	covariates	# $b$	# $O$	# $\Delta_j L$	total
$J$	$P$	$P$	$P \cdot (P + 1)/2$	$(J - 1) \cdot J/2 - 1$	
2	2	2	3	0	5
10	10	10	55	44	109

## Number of parameters to estimate

alternatives	covariates	#b	#O	# $\Delta_j L$	total
J	P	P	$P \cdot (P + 1)/2$	$(J - 1) \cdot J/2 - 1$	
2	2	2	3	0	5
10	10	10	55	44	109

## Log-likelihood

$$\log L(y) = \sum_{n,t,j} 1(y_{nt} = j) \overbrace{\Phi_{J-1}(-\Delta_j X_{nt} b \mid 0; \Delta_j (X_{nt} \Omega X_{nt}' + \Sigma) \Delta_j')}^{\Pr(y_{nt}=j)}$$

## 1 The multinomial probit model

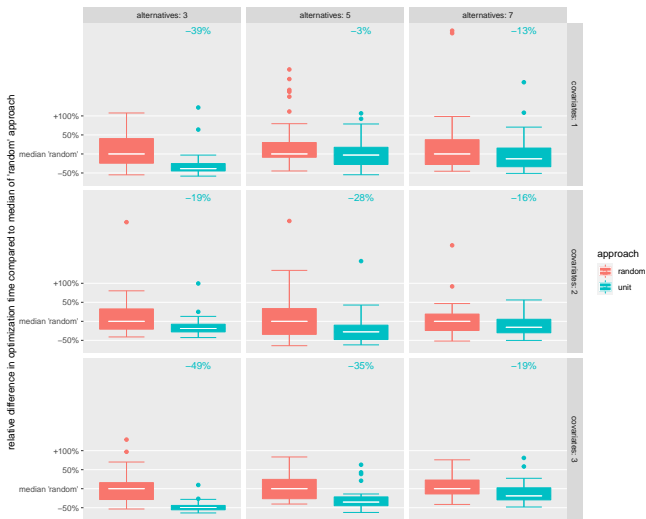
### 2 Initialization

- Unit
- Scaling
- Subsample
- Alternating optimization



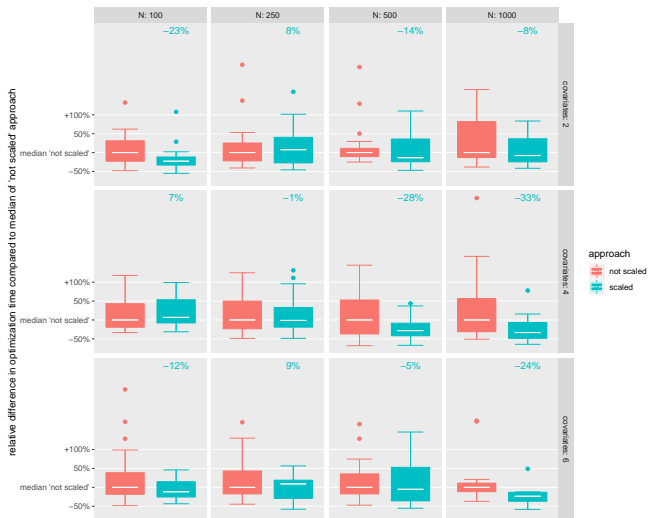
## Approach: Unit

- Idea: Choose initial parameters s.t.  $b = 0$ ,  $\Omega = I$  and  $\Sigma = I$
- Question: Better than random initialization?
- Simulation setting:
  - $N = 50$ ,  $T = 10$
  - random covariates: 1, 2, 3
  - alternatives: 3, 5, 7



## Approach: Scaling

- Idea: Different ranges of covariate values may hinder optimization (e.g. 112 minutes travel time, EUR 5 travel cost)
- Question: Does standardization improve optimization time?
- Simulation setting:
  - $T = 10$ , alternatives: 3
  - $N$ : 100, 250, 500, 1000
  - random covariates: 2, 4, 6
  - scale difference:  $U[1, 10]$
  - use same random initial guesses, but mind the scales



## Approach: Subsample

- Idea: Estimate the model on a subsample based on
  - random subsampling (rs)
  - k-means (km)
  - model-based clustering (mc)
- Simulation setting:
  - $T = 10$ , alternatives: 3, random covariates: 3
  - $N$ : 100, 250, 500, 1000
  - subsample proportion: 0.1, 0.25, 0.5



## Approach: Alternating optimization

- Idea: Alternating estimation of parameter groups
  - $b$  separately (ao\_b)
  - $b$  and variances separately (ao\_b\_var)
  - $b$ , variances and covariances separately (ao\_full)
- Simulation setting:
  - $N = 50, T = 10$
  - random covariates: 1, 2, 3
  - alternatives: 3, 5, 7





*Thank you!*

Please let me know:

- To what extent is initialization an issue for your models?
- How do you initialize?