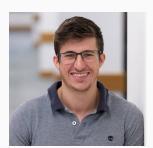
Filtering procedures for dynamic multinomial probit models

Augusto Fasano

University of Torino and Collegio Carlo Alberto

SISBAYES Workshop, 4 September 2025.

Joint work with:



Lorenzo Rimella

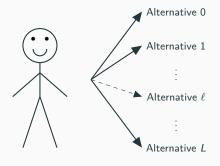


Giovanni Rebaudo

Introduction

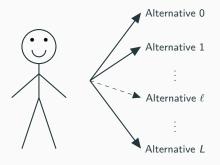
A brief overview on discrete-choice models

Discrete-choice models: A discrete choice model specifies the probability that a person chooses a particular alternative.



A brief overview on discrete-choice models

Discrete-choice models: A discrete choice model specifies the probability that a person chooses a particular alternative.



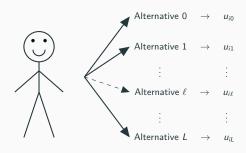
The probability of choosing a specific alternative is expressed as a function of observed variables that relate to:

- the person
- the alternatives

2

How is the choice made?

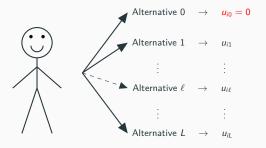
- Individual i obtains a utility $u_{i\ell}$ by choosing alternative ℓ , for $\ell=0,1,\ldots,L$.
- The behavior of the person is utility-maximizing: individual *i* chooses the alternative that provides the highest utility.



3

How is the choice made?

- Individual i obtains a utility $u_{i\ell}$ by choosing alternative ℓ , for $\ell = 0, 1, \dots, L$.
- The behavior of the person is utility-maximizing: individual *i* chooses the alternative that provides the highest utility.
- Only differences in utilities count. We set $u_{i0} = 0$ for identifiability.



How are utilities defined?

Calling $\mathbf{u}_i = (u_{i1}, \dots, u_{iL})^\mathsf{T}$, in the multinomial probit model (MNP) we have

$$\mathbf{u}_i = \mathbf{X}_i \boldsymbol{\beta} + \boldsymbol{\varepsilon}_i,$$

where

- X_i : $L \times p$ matrix which can contain
 - intercept terms,
 - covariates that vary across agents (e.g., age)
 - and/or covariates that vary across alternatives (e.g., prices)
- β : p-dimensional vector of parameters to be estimated.
- $\varepsilon_i \sim N_L(\mathbf{0}, \Sigma)$, where Σ is a covariance matrix satisfying appropriate constraints for identifiability (e.g., $\sigma_{11} = 1$ or $tr(\Sigma) = L$).
- Errors ε_{ij} , ε_{ik} , $j \neq k$, are not independent.

Calling $\mathbf{x}_{i\ell}^{\mathsf{T}}$ the ℓ -th row of \mathbf{X}_i , we thus have $u_{i\ell} = \mathbf{x}_{i\ell}^{\mathsf{T}} \boldsymbol{\beta} + \varepsilon_{i\ell}$.

4

Let us consider Σ fixed (for the moment) and call $\mathbf{x}_{i\ell}^{\mathsf{T}}$ the ℓ -th row of \mathbf{X}_i .

The likelihood of observation y_i is

$$\begin{aligned} \Pr[y_i = \ell \mid \beta] &= \Pr[u_{i\ell} > u_{ik} \ \forall k \neq \ell \mid \beta] \\ &= \Pr[\cap_{k \neq \ell} \{\mathbf{x}_{i\ell}^\mathsf{T} \boldsymbol{\beta} + \varepsilon_{i\ell} > \mathbf{x}_{ik}^\mathsf{T} \boldsymbol{\beta} + \varepsilon_{ik}\} \mid \beta] \\ &= \Pr[\cap_{k \neq \ell} \{\varepsilon_{ik} - \varepsilon_{i\ell} < (\mathbf{x}_{i\ell}^\mathsf{T} - \mathbf{x}_{ik})^\mathsf{T} \boldsymbol{\beta}\} \mid \beta] \\ &= \Pr[\mathbf{V}_{[-\ell]} \varepsilon_i < \mathbf{X}_{[i,-\ell]} \boldsymbol{\beta} \mid \beta], \end{aligned}$$

for appropriate matrices $\mathbf{V}_{[-\ell]}$ and $\mathbf{X}_{[i,-\ell]}$, with $\varepsilon_{i0}=0$ and $\mathbf{x}_{i0}=\mathbf{0}_p^\intercal$.

Let us consider Σ fixed (for the moment) and call $\mathbf{x}_{i\ell}^{\mathsf{T}}$ the ℓ -th row of \mathbf{X}_i .

The likelihood of observation y_i is

$$\begin{aligned} \Pr[y_i = \ell \mid \beta] &= \Pr[u_{i\ell} > u_{ik} \ \forall k \neq \ell \mid \beta] \\ &= \Pr[\cap_{k \neq \ell} \{\mathbf{x}_{i\ell}^\mathsf{T} \boldsymbol{\beta} + \varepsilon_{i\ell} > \mathbf{x}_{ik}^\mathsf{T} \boldsymbol{\beta} + \varepsilon_{ik}\} \mid \beta] \\ &= \Pr[\cap_{k \neq \ell} \{\varepsilon_{ik} - \varepsilon_{i\ell} < (\mathbf{x}_{i\ell}^\mathsf{T} - \mathbf{x}_{ik})^\mathsf{T} \boldsymbol{\beta}\} \mid \beta] \\ &= \Pr[\mathbf{V}_{[-\ell]} \varepsilon_i < \mathbf{X}_{[i,-\ell]} \boldsymbol{\beta} \mid \beta], \end{aligned}$$

for appropriate matrices $\mathbf{V}_{[-\ell]}$ and $\mathbf{X}_{[i,-\ell]}$, with $\varepsilon_{i0}=0$ and $\mathbf{x}_{i0}=\mathbf{0}_p^{\mathsf{T}}$.

MNP likelihood of a single observation

Thus, the likelihood of the *i*-th observation in the MNP model is

$$\Pr[y_i = \ell \mid \boldsymbol{\beta}] = \Phi_L(\boldsymbol{\mathsf{X}}_{[i,-\ell]}\boldsymbol{\beta};\boldsymbol{\mathsf{S}}_\ell), \qquad \boldsymbol{\mathsf{S}}_\ell = \boldsymbol{\mathsf{V}}_{[-\ell]}\boldsymbol{\Sigma}\boldsymbol{\mathsf{V}}_{[-\ell]}^\intercal,$$

where $\Phi_L(\mathbf{x}; \Sigma)$ denotes the cdf of a $N_L(\mathbf{0}, \Sigma)$ random v., evaluated at \mathbf{x} .

probit models

(Bayesian) Inference in the multinomial

MNP likelihood of a sample

The likelihood of a sample $\mathbf{y}_{1:n} = (y_1, \dots, y_n)$ is thus given by

$$p(\mathbf{y}_{1:n} \mid eta) = \prod_{i=1}^n \Phi_L(\mathbf{X}_{[i,-y_i]}eta; \mathbf{S}_{y_i}) = \Phi_{nL}(\mathbf{X}eta; oldsymbol{\Lambda}),$$

with

$$\mathbf{X} = [\mathbf{X}_{[1,-y_1]}^\intercal, \dots, \mathbf{X}_{[n,-y_n]}^\intercal]^\intercal, \qquad \boldsymbol{\Lambda} = \mathsf{block\text{-}diag}(\mathbf{S}_{y_1}, \dots, \mathbf{S}_{y_n})$$

How can we make inference on β ?

6

MNP likelihood of a sample

The likelihood of a sample $\mathbf{y}_{1:n} = (y_1, \dots, y_n)$ is thus given by

$$p(\mathbf{y}_{1:n} \mid \beta) = \prod_{i=1}^n \Phi_L(\mathbf{X}_{[i,-y_i]}\beta; \mathbf{S}_{y_i}) = \Phi_{nL}(\mathbf{X}\beta; \mathbf{\Lambda}),$$

with

$$\mathbf{X} = [\mathbf{X}_{[1,-y_1]}^\intercal, \dots, \mathbf{X}_{[n,-y_n]}^\intercal]^\intercal, \qquad \boldsymbol{\Lambda} = \mathsf{block\text{-}diag}(\mathbf{S}_{y_1}, \dots, \mathbf{S}_{y_n})$$

How can we make inference on β ?

We focus on the Bayesian framework, with multivariate Gaussian prior

$$eta \sim \mathcal{N}_{\scriptscriptstyle p}(\mu,\Omega)$$
 (usually $\mu=0$).

Thus

$$p(eta \mid \mathbf{y}_{1:n}) \propto d\mathcal{N}_p(eta; \mu, \Omega) \cdot \Phi_{nL}(\mathbf{X}eta; \mathbf{\Lambda}).$$

6

MNP likelihood of a sample

The likelihood of a sample $\mathbf{y}_{1:n} = (y_1, \dots, y_n)$ is thus given by

$$ho(\mathbf{y}_{1:n}\mideta)=\prod_{i=1}^n\Phi_L(\mathbf{X}_{[i,-y_i]}eta;\mathbf{S}_{y_i})=\Phi_{nL}(\mathbf{X}eta;oldsymbol{\Lambda}),$$

with

$$\mathbf{X} = [\mathbf{X}_{[1,-y_1]}^\intercal, \dots, \mathbf{X}_{[n,-y_n]}^\intercal]^\intercal, \qquad \boldsymbol{\Lambda} = \mathsf{block\text{-}diag}(\mathbf{S}_{y_1}, \dots, \mathbf{S}_{y_n})$$

How can we make inference on β ?

We focus on the Bayesian framework, with multivariate Gaussian prior

$$eta \sim \mathcal{N}_{\scriptscriptstyle p}(\mu,\Omega)$$
 (usually $\mu=0$).

Thus

$$p(\beta \mid \mathbf{y}_{1:n}) \propto d\mathcal{N}_p(\beta; \boldsymbol{\mu}, \boldsymbol{\Omega}) \cdot \Phi_{nL}(\mathbf{X}\beta; \boldsymbol{\Lambda}).$$

Is this the kernel of some known family of distributions?

Short recap: a $\mathrm{SUN}_{\rho,m}(\mu,\Omega,\Delta,\gamma,\Gamma)$ distributed random variable has density

$$p(eta) = d\mathcal{N}_p(eta; oldsymbol{\mu}, oldsymbol{\Omega}) rac{\Phi_m \left(oldsymbol{\gamma} + oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} \omega (eta - oldsymbol{\xi}); oldsymbol{\Gamma} - oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} oldsymbol{\Delta}
ight)}{\Phi_m (oldsymbol{\gamma}; oldsymbol{\Gamma})},$$

where $\bar{\Omega}$ and ω correlation and scale matrices associated to $\Omega = \omega \bar{\Omega} \omega$; Moments may be quite involved, but we can develop an **i.i.d.** sampler.

Important fact: SUN additive representation

If $eta \sim \mathrm{SUN}_{p,m}$, then it can be characterized probabilistically as a linear combination of:

- a p-variate Gaussian term;
- an *m*-variate truncated Gaussian component.

Short recap: a $\mathrm{SUN}_{\rho,m}(\mu,\Omega,\Delta,\gamma,\Gamma)$ distributed random variable has density

$$p(eta) = d\mathcal{N}_p(eta; oldsymbol{\mu}, oldsymbol{\Omega}) rac{\Phi_m \left(\gamma + oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} oldsymbol{\omega} (eta - oldsymbol{\xi}); oldsymbol{\Gamma} - oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} oldsymbol{\Delta}
ight)}{\Phi_m(\gamma; oldsymbol{\Gamma})}$$

where $\bar{\Omega}$ and ω correlation and scale matrices associated to $\Omega = \omega \bar{\Omega} \omega$; Moments may be quite involved, but we can develop an **i.i.d.** sampler.

Important fact: SUN additive representation

If $eta \sim \mathrm{SUN}_{p,m}$, then it can be characterized probabilistically as a linear combination of:

- a p-variate Gaussian term;
- an *m*-variate truncated Gaussian component.

In the MNP, we have $p(\beta \mid \mathbf{y}_{1:n}) \propto d\mathcal{N}_p(\beta; \mu, \Omega) \cdot \Phi_{nL}(\mathbf{X}\beta; \Lambda)$ $\implies p(\beta \mid \mathbf{y}_{1:n})$ is the kernel of a $\mathrm{SUN}_{p,nL}$ (F., Durante, JMLR, 2022).

Short recap: a $\mathrm{SUN}_{\rho,m}(\mu,\Omega,\Delta,\gamma,\Gamma)$ distributed random variable has density

$$p(eta) = d\mathcal{N}_p(eta; oldsymbol{\mu}, oldsymbol{\Omega}) rac{\Phi_m \left(oldsymbol{\gamma} + oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} \omega (eta - oldsymbol{\xi}); \Gamma - oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} oldsymbol{\Delta}
ight)}{\Phi_m (oldsymbol{\gamma}; \Gamma)},$$

where $\bar{\Omega}$ and ω correlation and scale matrices associated to $\Omega = \omega \bar{\Omega} \omega$; Moments may be quite involved, but we can develop an **i.i.d.** sampler.

Important fact: SUN additive representation

If $eta \sim \mathrm{SUN}_{p,m}$, then it can be characterized probabilistically as a linear combination of:

- a p-variate Gaussian term;
- an *m*-variate truncated Gaussian component.

In the MNP, we have $p(\beta \mid \mathbf{y}_{1:n}) \propto d\mathcal{N}_p(\beta; \boldsymbol{\mu}, \boldsymbol{\Omega}) \cdot \Phi_{nL}(\mathbf{X}\beta; \boldsymbol{\Lambda})$

- $\implies p(\beta \mid \mathbf{y}_{1:n})$ is the kernel of a $SUN_{p,nL}$ (F., Durante, JMLR, 2022).
- \implies we get an i.i.d. sampler for $p(\beta \mid \mathbf{y}_{1:n})$ well-suited for high-dimensional scenarios, but it becomes infeasible as nL gets larger.

Short recap: a $\mathrm{SUN}_{\rho,m}(\mu,\Omega,\Delta,\gamma,\Gamma)$ distributed random variable has density

$$p(eta) = d\mathcal{N}_p(eta; oldsymbol{\mu}, oldsymbol{\Omega}) rac{\Phi_m \left(\gamma + oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} oldsymbol{\omega} (eta - oldsymbol{\xi}); oldsymbol{\Gamma} - oldsymbol{\Delta}^\intercal ar{oldsymbol{\Omega}}^{-1} oldsymbol{\Delta}
ight)}{\Phi_m (\gamma; oldsymbol{\Gamma})},$$

where $\bar{\Omega}$ and ω correlation and scale matrices associated to $\Omega = \omega \bar{\Omega} \omega$; Moments may be quite involved, but we can develop an **i.i.d.** sampler.

Important fact: SUN additive representation

If $\beta \sim \mathrm{SUN}_{p,m}$, then it can be characterized probabilistically as a linear combination of:

- a p-variate Gaussian term;
- an *m*-variate truncated Gaussian component.

In the MNP, we have $p(\beta \mid \mathbf{y}_{1:n}) \propto d\mathcal{N}_p(\beta; \mu, \Omega) \cdot \Phi_{nL}(\mathbf{X}\beta; \Lambda)$

- $\implies p(\beta \mid \mathbf{y}_{1:n})$ is the kernel of a $SUN_{p,nL}$ (F., Durante, JMLR, 2022).
- \implies we get an i.i.d. sampler for $p(\beta \mid \mathbf{y}_{1:n})$ well-suited for high-dimensional scenarios, but it becomes infeasible as nL gets larger.
 - ⇒ Variational Bayes can help to overcome this issues.

The dynamic multinomial probit model

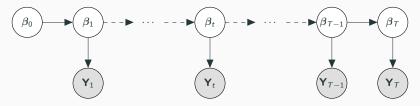
The dynamic multinomial probit model

What can we say about **dynamic versions** of the MNP model that account for time-dependent observations?

The dynamic multinomial probit model

What can we say about **dynamic versions** of the MNP model that account for time-dependent observations?

Dynamic multinomial probit model:



 $\mathbf{Y}_t = (y_{t,1}, \dots, y_{t,n_t})$ represents the n_t categorical observations sampled at time t.

with $\beta_0 \sim \mathcal{N}_p \left(\mathbf{a}_0, \mathbf{P}_0 \right)$ independent of $\boldsymbol{\eta}_t \stackrel{\mathsf{iid}}{\sim} \mathcal{N}_p \left(\mathbf{0}, \mathbf{W} \right)$.

R

Methodological question:

Can we develop **online procedures** for the **filtering** distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$?

Methodological question:

Can we develop **online procedures** for the **filtering** distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$?

Rimella, F. and Rebaudo (2025+)

Calling $N_t = \sum_{s=1}^t n_s$ the total number of observations up to time t, under the dynamic probit model:

- 1. the filtering distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$ is SUN_{p,N_tL} ;
- 2. the state predictive distribution $p(\beta_{t+1} \mid \mathbf{Y}_{1:t})$ is SUN_{p,N_tL} ;

Methodological question:

Can we develop **online procedures** for the **filtering** distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$?

Rimella, F. and Rebaudo (2025+)

Calling $N_t = \sum_{s=1}^t n_s$ the total number of observations up to time t, under the dynamic probit model:

- 1. the filtering distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$ is SUN_{p,N_tL} ;
- 2. the state predictive distribution $p(\beta_{t+1} \mid \mathbf{Y}_{1:t})$ is SUN_{p,N_tL} ;

l.i.d. sampling in principle feasible thanks to the $\mathop{\rm SUN}\nolimits$ additive representation. Need to linearly combine:

- a sample from a p-variate Gaussian;
- a sample from a N_tL -variate truncated Gaussian.

The truncated multivariate component is the computationally-impractical part.

Methodological question:

Can we develop **online procedures** for the **filtering** distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$?

Rimella, F. and Rebaudo (2025+)

Calling $N_t = \sum_{s=1}^t n_s$ the total number of observations up to time t, under the dynamic probit model:

- 1. the filtering distribution $p(\beta_t \mid \mathbf{Y}_{1:t})$ is SUN_{p,N_tL} ;
- 2. the state predictive distribution $p(\beta_{t+1} \mid \mathbf{Y}_{1:t})$ is SUN_{p,N_tL} ;

l.i.d. sampling in principle feasible thanks to the $\mathop{\rm SUN}\nolimits$ additive representation. Need to linearly combine:

- a sample from a p-variate Gaussian;
- a sample from a N_tL -variate truncated Gaussian.

The truncated multivariate component is the computationally-impractical part.

Can we get a procedure that scales linearly in t?

Computational methods

Basic idea: at each time t, obtain a sample from the **target distribution** $p(\beta_{1:t} \mid \mathbf{Y}_{1:t})$, exploiting the samples obtained at previous iterations t-1.

Basic idea: at each time t, obtain a sample from the **target distribution** $p(\beta_{1:t} \mid \mathbf{Y}_{1:t})$, exploiting the samples obtained at previous iterations t-1.

The performance of particle filters relies on the **proposal** $\pi(\beta_{t|t} \mid \beta_{1:t-1|1:t-1}^{(r)}, \mathbf{Y}_{1:t})$ and the form of the **resampling weights** $w_t^{(r)} = w_t(\bar{\beta}_{1:t|t}^{(r)}), r = 1, \dots, R$.

Basic idea: at each time t, obtain a sample from the **target distribution** $p(\beta_{1:t} \mid \mathbf{Y}_{1:t})$, exploiting the samples obtained at previous iterations t-1.

The performance of particle filters relies on the **proposal** $\pi(\beta_{t|t} \mid \beta_{1:t-1|1:t-1}^{(r)}, \mathbf{Y}_{1:t})$ and the form of the **resampling weights** $w_t^{(r)} = w_t(\bar{\beta}_{1:t|t}^{(r)}), r = 1, \dots, R$.

In the simplest case (bootstrap particle filter):

- $\pi(\beta_t \mid \beta_{1:t-1}, \mathbf{Y}_{1:t}) = p(\beta_t \mid \beta_{t-1})$
- $\qquad \qquad \mathbf{w}_t^{(r)} \propto p(\mathbf{Y}_t \mid \bar{\beta}_t^{(r)}) = \Phi_L(\mathbf{X}_{[y_t,1]} \bar{\beta}_{t|t}^{(r)}; \mathbf{S}_{y_{t,1}}) \cdots \Phi_L(\mathbf{X}_{[y_{t,n_t}]} \bar{\beta}_{t|t}^{(r)}; \mathbf{S}_{y_{t,n_t}})$

At each step t, we would need to compute $R \times n_t$ L-dimensional Gaussian cdfs.

Basic idea: at each time t, obtain a sample from the **target distribution** $p(\beta_{1:t} \mid \mathbf{Y}_{1:t})$, exploiting the samples obtained at previous iterations t-1.

The performance of particle filters relies on the **proposal** $\pi(\beta_{t|t} \mid \beta_{1:t-1|1:t-1}^{(r)}, \mathbf{Y}_{1:t})$ and the form of the **resampling weights** $w_t^{(r)} = w_t(\bar{\beta}_{1:t|t}^{(r)}), r = 1, \dots, R$.

In the simplest case (bootstrap particle filter):

$$\pi(\beta_t \mid \beta_{1:t-1}, \mathbf{Y}_{1:t}) = p(\beta_t \mid \beta_{t-1})$$

$$w_t^{(r)} \propto p(\mathbf{Y}_t \mid \bar{\beta}_t^{(r)}) = \Phi_L(\mathbf{X}_{[y_{t-1}]} \bar{\beta}_{t|t}^{(r)}; \mathbf{S}_{y_{t-1}}) \cdots \Phi_L(\mathbf{X}_{[y_{t-n_t}]} \bar{\beta}_{t|t}^{(r)}; \mathbf{S}_{y_{t,n_t}})$$

At each step t, we would need to compute $R \times n_t$ L-dimensional Gaussian cdfs.

 \implies computationally demanding.

Basic idea: at each time t, obtain a sample from the **target distribution** $p(\beta_{1:t} \mid \mathbf{Y}_{1:t})$, exploiting the samples obtained at previous iterations t-1.

The performance of particle filters relies on the **proposal** $\pi(\beta_{t|t} \mid \beta_{1:t-1|1:t-1}^{(r)}, \mathbf{Y}_{1:t})$ and the form of the **resampling weights** $w_t^{(r)} = w_t(\bar{\beta}_{1:t|t}^{(r)}), r = 1, \dots, R$.

In the simplest case (bootstrap particle filter):

- $\quad \boldsymbol{\pi}(\boldsymbol{\beta}_t \mid \boldsymbol{\beta}_{1:t-1}, \mathbf{Y}_{1:t}) = \boldsymbol{p}(\boldsymbol{\beta}_t \mid \boldsymbol{\beta}_{t-1})$
- $\qquad \qquad \mathbf{w}_t^{(r)} \propto p(\mathbf{Y}_t \mid \bar{\beta}_t^{(r)}) = \Phi_L(\mathbf{X}_{[y_t,1]} \bar{\beta}_{t|t}^{(r)}; \mathbf{S}_{y_{t,1}}) \cdots \Phi_L(\mathbf{X}_{[y_{t,n_t}]} \bar{\beta}_{t|t}^{(r)}; \mathbf{S}_{y_{t,n_t}})$

At each step t, we would need to compute $R \times n_t$ L-dimensional Gaussian cdfs.

 \implies computationally demanding.

Other choices for the proposal would be possible (optimal proposal as in Doucet et al., 2000), but the weigths would still have the same problem.

Basic idea: use sequential **simple** SUN approximations $h_t(\beta_t)$ of the **filtering distribution** $p(\beta_t \mid \mathbf{Y}_{1:t})$, by sequentially approximating the filtering distribution of the previous time with a Gaussian density.

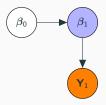


At time
$$t = 0$$
: $p(\beta_0) = d\mathcal{N}_p(\beta_0; \mathbf{a}_0, \mathbf{P}_0)$.



At time t = 0: $p(\beta_0) = d\mathcal{N}_p(\beta_0; \mathbf{a}_0, \mathbf{P}_0)$.

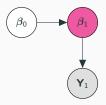
This gives a prior distribution for β_1 : $p(\beta_1) = d\mathcal{N}_p(\beta_1; \mathbf{Ga}_0, \mathbf{GP}_0\mathbf{G}^\mathsf{T} + \mathbf{W})$.



At time t = 0: $p(\beta_0) = d\mathcal{N}_p(\beta_0; \mathbf{a}_0, \mathbf{P}_0)$.

This gives a prior distribution for β_1 : $p(\beta_1) = d\mathcal{N}_p(\beta_1; Ga_0, GP_0G^T + W)$.

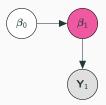
At time t=1: $p(\beta_1 \mid \mathbf{Y}_1) \propto p(\beta_1) p(\mathbf{Y}_1 \mid \beta_1)$



At time t = 0: $p(\beta_0) = d\mathcal{N}_p(\beta_0; \mathbf{a}_0, \mathbf{P}_0)$.

This gives a prior distribution for β_1 : $p(\beta_1) = d\mathcal{N}_p(\beta_1; Ga_0, GP_0G^T + W)$.

At time t = 1: $p(\beta_1 \mid \mathbf{Y}_1) \propto p(\beta_1) \Phi_{n_1 L}(\mathbf{X}_{[\mathbf{Y}_1]} \beta_1; \mathbf{\Lambda}_{[\mathbf{Y}_1]}) = d \text{SUN}_{p, n_1 L}$

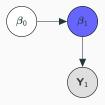


At time t = 0: $p(\beta_0) = d\mathcal{N}_p(\beta_0; \mathbf{a}_0, \mathbf{P}_0)$.

This gives a prior distribution for β_1 : $p(\beta_1) = d\mathcal{N}_p(\beta_1; Ga_0, GP_0G^T + W)$.

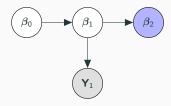
At time t=1: $p(\beta_1 \mid \mathbf{Y}_1) \propto p(\beta_1) \Phi_{n_1 L}(\mathbf{X}_{[\mathbf{Y}_1]} \beta_1; \mathbf{\Lambda}_{[\mathbf{Y}_1]}) = d \mathrm{SUN}_{\rho, n_1 L}$

Idea: Approximate now $p(\beta_1 | Y_1)$ with a Gaussian density $q_1(\beta_1)$ and repeat the process at next time.



At time t = 2:

Approximate $p(\beta_1 \mid \mathbf{Y}_1)$ with $q_1(\beta_1) = d\mathcal{N}_p(\beta_1; \mu_{1|1}, \Omega_{1|1})$ matching the first two moments.

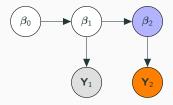


At time t = 2:

Approximate $p(\beta_1 \mid \mathbf{Y}_1)$ with $q_1(\beta_1) = d\mathcal{N}_p(\beta_1; \mu_{1|1}, \Omega_{1|1})$ matching the first two moments.

This induces a **Gaussian predictive** distribution for β_2 :

$$q_1(eta_2) = d\mathcal{N}_{
ho}(eta_2; \mathbf{G}\mu_{1|1}, \mathbf{G}\Omega_{1|1}\mathbf{G}^\intercal + \mathbf{W}) pprox
ho(eta_2 \mid \mathbf{Y_1}).$$



At time t = 2:

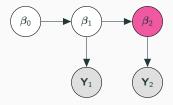
Approximate $p(\beta_1 \mid \mathbf{Y}_1)$ with $q_1(\beta_1) = d\mathcal{N}_p(\beta_1; \mu_{1|1}, \Omega_{1|1})$ matching the first two moments.

This induces a **Gaussian predictive** distribution for β_2 :

$$q_1(eta_2) = d\mathcal{N}_{
ho}(eta_2; \mathbf{G}\mu_{1|1}, \mathbf{G}\Omega_{1|1}\mathbf{G}^{\mathsf{T}} + \mathbf{W}) pprox
ho(eta_2 \mid \mathbf{Y}_1).$$

Plugging this in $p(\beta_2 \mid \mathbf{Y}_{1:2}) \propto p(\beta_2 \mid \mathbf{Y}_1) p(\mathbf{Y}_2 \mid \beta_2)$, we get the approximation

$$p(\beta_2 \mid \mathbf{Y}_{1:2}) \approx h_2(\beta_2) \propto q_1(\beta_2) p(\mathbf{Y}_2 \mid \beta_2)$$



At time t = 2:

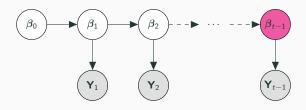
Approximate $p(\beta_1 \mid \mathbf{Y}_1)$ with $q_1(\beta_1) = d\mathcal{N}_p(\beta_1; \mu_{1|1}, \Omega_{1|1})$ matching the first two moments.

This induces a **Gaussian predictive** distribution for β_2 :

$$q_1(eta_2) = d\mathcal{N}_p(eta_2; \mathbf{G}\mu_{1|1}, \mathbf{G}\Omega_{1|1}\mathbf{G}^\intercal + \mathbf{W}) pprox p(eta_2 \mid \mathbf{Y_1}).$$

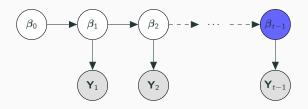
Plugging this in $p(\beta_2 \mid \mathbf{Y}_{1:2}) \propto p(\beta_2 \mid \mathbf{Y}_1) p(\mathbf{Y}_2 \mid \beta_2)$, we get the approximation

$$p(\beta_2 \mid \mathbf{Y}_{1:2}) \approx h_2(\beta_2) \propto q_1(\beta_2) p(\mathbf{Y}_2 \mid \beta_2) = dSUN_{p,n_2L}.$$



At generic time t:

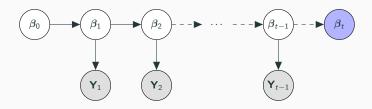
From previous time t-1, we have the **SUN** filtering approximation $h_{t-1}(\beta_{t-1})$.



At generic time t:

From previous time t-1, we have the **SUN** filtering approximation $h_{t-1}(\beta_{t-1})$.

Then, take the Gaussian approximation $q_{t-1}(\beta_t) \approx p(\beta_{t-1} \mid \mathbf{Y}_{1:t-1})$, matching the first two moments of $h_{t-1}(\beta_{t-1})$.

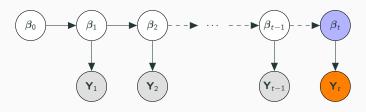


At generic time t:

From previous time t-1, we have the **SUN** filtering approximation $h_{t-1}(\beta_{t-1})$.

Then, take the Gaussian approximation $q_{t-1}(\beta_t) \approx p(\beta_{t-1} \mid \mathbf{Y}_{1:t-1})$, matching the first two moments of $h_{t-1}(\beta_{t-1})$.

This induces a Gaussian predictive distribution for β_t : $q_{t-1}(\beta_t) \approx p(\beta_t \mid \mathbf{Y}_{1:t-1})$.



At generic time *t*:

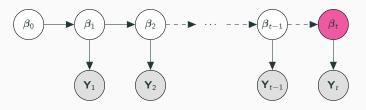
From previous time t-1, we have the **SUN** filtering approximation $h_{t-1}(\beta_{t-1})$.

Then, take the **Gaussian** approximation $q_{t-1}(\beta_t) \approx p(\beta_{t-1} \mid \mathbf{Y}_{1:t-1})$, matching the first two moments of $h_{t-1}(\beta_{t-1})$.

This induces a **Gaussian predictive** distribution for β_t : $q_{t-1}(\beta_t) \approx p(\beta_t \mid \mathbf{Y}_{1:t-1})$.

Plugging this in $p(\beta_t \mid \mathbf{Y}_{1:t}) \propto p(\beta_t \mid \mathbf{Y}_{1:t-1}) p(\mathbf{Y}_t \mid \beta_t)$, we get

$$p(\beta_t \mid \mathbf{Y}_{1:t}) \approx h_t(\beta_t) \propto q_{t-1}(\beta_t) p(\mathbf{Y}_t \mid \beta_t)$$



At generic time *t*:

From previous time t-1, we have the **SUN** filtering approximation $h_{t-1}(\beta_{t-1})$.

Then, take the **Gaussian** approximation $q_{t-1}(\beta_t) \approx p(\beta_{t-1} \mid \mathbf{Y}_{1:t-1})$, matching the first two moments of $h_{t-1}(\beta_{t-1})$.

This induces a Gaussian predictive distribution for β_t : $q_{t-1}(\beta_t) \approx p(\beta_t \mid \mathbf{Y}_{1:t-1})$.

Plugging this in $p(\beta_t \mid \mathbf{Y}_{1:t}) \propto p(\beta_t \mid \mathbf{Y}_{1:t-1}) p(\mathbf{Y}_t \mid \boldsymbol{\beta}_t)$, we get

$$p(\beta_t \mid \mathbf{Y}_{1:t}) \approx h_t(\beta_t) \propto q_{t-1}(\beta_t) p(\mathbf{Y}_t \mid \beta_t) = dSUN_{\rho, n_t L}$$

Wrap up of methods

Computational bottlenecks of the various methods seen so far:

- i.i.d. sampler: sampling from $N_t \times L$ -variate truncated Gaussian, $N_t = \sum_{s=1}^t n_s$;
- particle filter: computation of $R \times n_t$ cdfs of L-variate Gaussians at each time;
- assumed density filtering (ADF): sampling from an $n_t \times L$ -variate truncated Gaussian at each time t (arising from the $SUN_{p,n_t \times L}$ approximated filtering distribution of time t).

Wrap up of methods

Computational bottlenecks of the various methods seen so far:

- i.i.d. sampler: sampling from $N_t \times L$ -variate truncated Gaussian, $N_t = \sum_{s=1}^t n_s$;
- particle filter: computation of $R \times n_t$ cdfs of L-variate Gaussians at each time;
- assumed density filtering (ADF): sampling from an $n_t \times L$ -variate truncated Gaussian at each time t (arising from the $\mathrm{SUN}_{p,n_t \times L}$ approximated filtering distribution of time t).

Can we get something more efficient? Can we use the spirit of ADF but avoid the $n_t \times L$ -variate truncated Gaussian?

Option 3: (Sequential) Expectation Propagation

Basic idea: see the filtering distribution

$$p(\beta_t \mid \mathbf{Y}_{1:t}) \propto p(\beta_t \mid \mathbf{Y}_{1:t-1}) \prod_{i=1}^{n_t} p(y_{t,i} \mid \beta_t)$$

as the posterior distribution in a model where

- $p(\beta_t \mid \mathbf{Y}_{1:t-1})$ is the prior,
- $p(y_{t,i} \mid \beta_t)$ is the likelihood of observation $y_{i,t}$, $i = 1, \ldots, n_t$.

Approximate this posterior via expectation propagation (EP) with a Gaussian approximation $q_t(\beta_t) = d\mathcal{N}_p(\beta_t; \mu_{t|t}, \Omega_{t|t})$.

Option 3: (Sequential) Expectation Propagation

In EP, we approximate $p(\beta_t \mid \mathbf{Y}_{1:t}) \propto p(\beta_t \mid \mathbf{Y}_{1:t-1}) \prod_{i=1}^{n_t} p(y_{t,i} \mid \boldsymbol{\beta}_t)$ with

$$q_t(eta_t) \propto q_{t-1}(eta_t) \prod_{i=1}^{n_t} q_{t,i}(eta_t),$$
 where

- $\quad \bullet \quad q_{t-1}(\beta_t) = d\mathcal{N}_{\rho}(\beta_t; \mathbf{G}\boldsymbol{\mu}_{t-1|t-1}, \mathbf{G}\boldsymbol{\Omega}_{t-1|t-1}\mathbf{G}^\mathsf{T} + \mathbf{W}) \approx \rho(\beta_t \mid \mathbf{Y}_{1:t-1}),$
- $q_{t,i}(\boldsymbol{\beta}_t) \propto \exp\{-\frac{1}{2}\beta_t^\intercal \mathbf{Q}_{t,i}\beta_t + \beta_t^\intercal \mathbf{r}_{t,i}\} \approx p(y_{t,i} \mid \boldsymbol{\beta}_t) \text{ for } i = 1, \dots, n_t.$

Since we have Gaussian prior and Gaussian likelihood

$$q_t(\boldsymbol{\beta}_t) = d\mathcal{N}_p(\boldsymbol{\beta}_t; \boldsymbol{\mu}_{t|t}, \boldsymbol{\Omega}_{t|t}),$$

where $\mu_{t|t}$ and $\Omega_{t|t}$ are determined by $\mathbf{r}_{t,i}$ and $\mathbf{Q}_{t,i}$, $i=1,\ldots,n_t$.

Option 3: (Sequential) Expectation Propagation

In EP, we approximate $p(\beta_t \mid \mathbf{Y}_{1:t}) \propto p(\beta_t \mid \mathbf{Y}_{1:t-1}) \prod_{i=1}^{n_t} p(y_{t,i} \mid \beta_t)$ with

$$q_t(eta_t) \propto q_{t-1}(eta_t) \prod_{i=1}^{n_t} q_{t,i}(eta_t),$$
 where

- $\quad \bullet \quad q_{t-1}(\beta_t) = d\mathcal{N}_{\rho}(\beta_t; \mathbf{G}\boldsymbol{\mu}_{t-1|t-1}, \mathbf{G}\boldsymbol{\Omega}_{t-1|t-1}\mathbf{G}^\mathsf{T} + \mathbf{W}) \approx \rho(\beta_t \mid \mathbf{Y}_{1:t-1}),$
- $q_{t,i}(\boldsymbol{\beta}_t) \propto \exp\{-\frac{1}{2}\beta_t^\mathsf{T} \mathbf{Q}_{t,i}\beta_t + \beta_t^\mathsf{T} \mathbf{r}_{t,i}\} \approx p(y_{t,i} \mid \boldsymbol{\beta}_t) \text{ for } i = 1, \dots, n_t.$

Since we have Gaussian prior and Gaussian likelihood

$$q_t(\boldsymbol{\beta}_t) = d\mathcal{N}_p(\boldsymbol{\beta}_t; \boldsymbol{\mu}_{t|t}, \boldsymbol{\Omega}_{t|t}),$$

where $\mu_{t|t}$ and $\Omega_{t|t}$ are determined by $\mathbf{r}_{t,i}$ and $\mathbf{Q}_{t,i}$, $i=1,\ldots,n_t$.

The quantities $\mathbf{r}_{t,i}$ and $\mathbf{Q}_{t,i}$, $i=1,\ldots,n_t$ are fixed in an "optimal" way via EP moment matching conditions, available in closed form.

Marginal likelihood via EP

EP can also be used to get an approximation of the marginal likelihood

$$p(\mathbf{Y}_{1:T}) = p(\mathbf{Y}_1)p(\mathbf{Y}_2 \mid \mathbf{Y}_1)p(\mathbf{Y}_3 \mid \mathbf{Y}_{1:2}) \cdots p(\mathbf{Y}_T \mid \mathbf{Y}_{1:T-1})$$

since at each time it can give an approximation of $p(\mathbf{Y}_t \mid \mathbf{Y}_{1:t-1})$.

Expectation-maximization (EM) can be used to **estimate** the desired **hyper-parameters** (e.g., Σ) by maximizing the marginal likelihood.

Experiments and results

We checked the ability of the proposed EM procedure to maximize the log-marginal likelihood and obtain meaningful estimates for Σ .

We ran the following experiment generating synthetic data as follows:

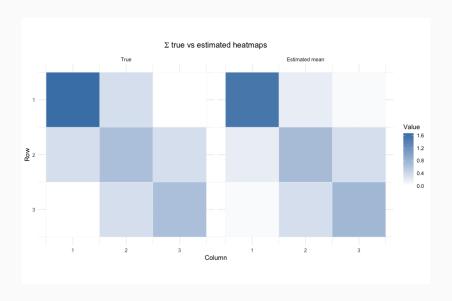
- we took L=3 and T=10,
- 15 individuals on average at each time step,
- 2 individual-specific covariates + 3 choice-specific covariates (+ choice-specific intercepts).

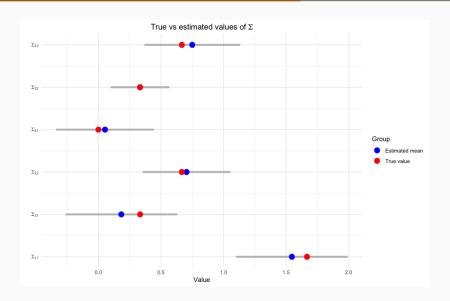
We checked the ability of the proposed EM procedure to maximize the log-marginal likelihood and obtain meaningful estimates for Σ .

We ran the following experiment generating synthetic data as follows:

- we took L=3 and T=10,
- 15 individuals on average at each time step,
- 2 individual-specific covariates + 3 choice-specific covariates
 (+ choice-specific intercepts).
- Keeping the above covariates fixed, we generated 50 datasets from the model:
 - each time generating a different trajectory of the parameters $\beta_{0:T}$,
 - then, given the generated $\beta_{0:T}$, we generated observations $\mathbf{Y}_{1:T}$.

Is the EM effective on average in the estimation of Σ ?





Filter experiments

We then compared the performance of the algorithms in a setting where:

- T = 20,
- $n_t \sim \text{Pois}(10)$ for each t,
- L = 3.
- there are 2 individual-specific characteristics,
- there are 3 choice-specific characteristics.

That is,

$$u_{t,i,\ell} = \alpha_{t,0,\ell} + \boldsymbol{\xi}_{t,i}^{\mathsf{T}} \boldsymbol{\alpha}_{t,\ell} + \boldsymbol{\zeta}_{t,i,\ell}^{\mathsf{T}} \boldsymbol{\gamma}_t + \varepsilon_{t,i,\ell}, \quad \ell = 1,\ldots,3,$$

where $\xi_{t,i}$ has dimension 2 and $\zeta_{t,i,\ell}$ has dimension 3.

The vector of parameters is then

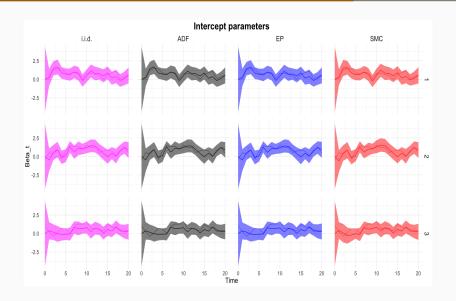
$$\boldsymbol{\beta}_t = (\alpha_{t,0,1}, \alpha_{t,0,2}, \alpha_{t,0,3}, \boldsymbol{\alpha}_{t,1}^\mathsf{T}, \boldsymbol{\alpha}_{t,2}^\mathsf{T}, \boldsymbol{\alpha}_{t,3}^\mathsf{T}, \gamma_t^\mathsf{T})^\mathsf{T} \in \mathbb{R}^{12}.$$

Filter experiments: running times

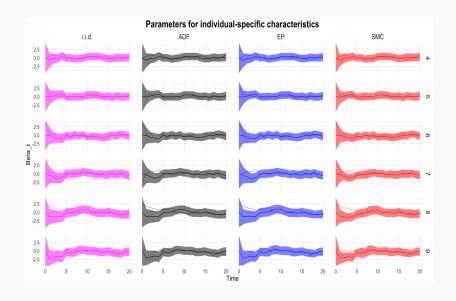
Method	i.i.d.	ADF	EP	SMC
Time (in seconds)	66,766.10	7.05	17.23	528.60

Table 1: Running times based on 5,000 samples.

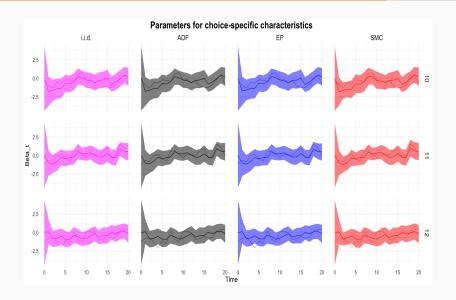
Filter experiments: intercepts parameters



Filter experiments: individual-specific characteristics parameters



Filter experiments: choice-specific characteristics parameters



Filter experiments: higher dimension

If we increase the dimension of the problem, taking

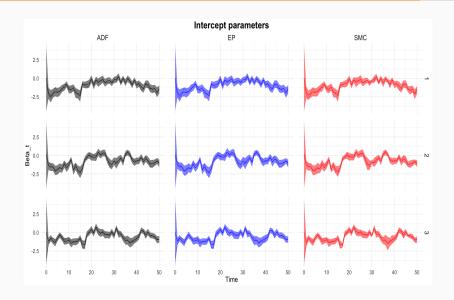
- T = 50,
- $n_t \sim \text{Pois}(50)$ for each t,
- the rest as before:
 - L = 3.
 - 2 individual-specific characteristics,
 - 3 choice-specific characteristics.

we obtain the following running times

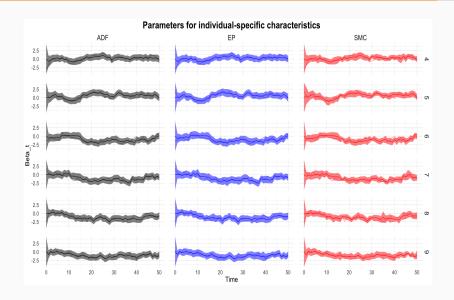
Method	i.i.d.	ADF	EP	SMC
Time (in seconds)	NA	599.98	199.69	6323.35

Table 2: Running times based on 5,000 samples.

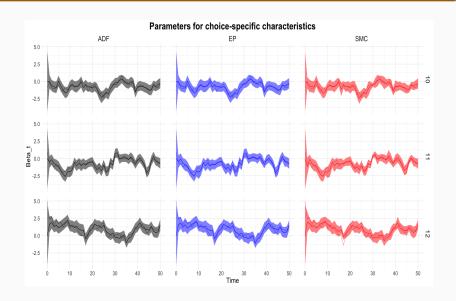
Filter experiments: intercepts parameters



Filter experiments: individual-specific characteristics parameters



Filter experiments: choice-specific characteristics parameters



Conclusions

- We have considered multiple filtering methods to make inference in dynamic discrete-choice models based on multinomial probit likelihoods.
- The i.i.d. sampler seems impractical already in moderate dimensions.
- Sequential Monte Carlo (SMC) procedures do not seem an appropriate solution due to the need to compute the MNP likelihood multiple times.
- The assumed density filtering (ADF) procedure is efficient in scenarios where n_t × L never exceeds a few hundred.
- In higher dimensions, sequential expectation propagation (EP) is the most efficient method, with accuracy comparable to the one of ADF.
- An expectation-maximization (EM) can be used to estimate the hyperparameters, like, e.g., the covariance matrix of the error terms Σ .
- Ongoing work on an Expedia dataset about bookings of properties across time: arxiv preprint coming in the coming months!

Thank You!