

Designing Graph Convolutional Neural Networks for Discrete Choice with Network Effects

Daniel F. Villarraga^a, Ricardo A. Daziano^a

^a*School of Civil and Environmental Engineering, Cornell University, 220 Hollister Hall, Ithaca, NY 14853, USA*

Abstract

We introduce a novel model architecture that incorporates network effects into discrete choice problems, achieving higher predictive performance than standard discrete choice models while offering greater interpretability than general-purpose flexible model classes. Econometric discrete choice models aid in studying individual decision-making, where agents select the option with the highest reward from a discrete set of alternatives. Intuitively, the utility an individual derives from a particular choice depends on their personal preferences and characteristics, the attributes of the alternative, and the value their peers assign to that alternative or their previous choices. However, most applications ignore peer influence, and models that do consider peer or network effects often lack the flexibility and predictive performance of recently developed approaches to discrete choice, such as deep learning. We propose a novel graph convolutional neural network architecture to model network effects in discrete choices, achieving higher predictive performance than standard discrete choice models while retaining the interpretability necessary for inference—a quality often lacking in general-purpose deep learning architectures. We evaluate our architecture using revealed commuting choice data, extended with travel times and trip costs for each travel mode for work-related trips in New York City, as well as 2016 U.S. election data aggregated by county, to test its performance on datasets with highly imbalanced classes. Given the interpretability of our models, we can estimate relevant economic metrics, such as the value of travel time savings in New York City. Finally, we compare the predictive performance and behavioral insights from our architecture to those derived from traditional discrete choice and general-purpose deep learning models.

Keywords: Network effects, Value of Travel Time Savings, Mode choice, Interpretable Deep Learning, Graph Neural Networks

Email addresses: dv275@cornell.edu (Daniel F. Villarraga), daziano@cornell.edu (Ricardo A. Daziano)

1. Introduction

Individuals tend to face choice situations where they have to select among a discrete set of alternatives, such as a mode of transportation for commuting to work, a candidate in political elections, or a new laptop to purchase. Econometric discrete choice models (DCMs) study these types of choice situations, assuming that individuals seek to maximize utility by selecting the most rewarding option. Typically, for standard discrete choice models, this choice depends on both the characteristics and preferences of the individual, and the attributes of each available option. Standard discrete choice models are interpretable and are used to extract microeconomic information, such as marginal rates of substitution (including willingness to pay for quality improvements), expected market shares, and probability marginal effects and choice elasticities with respect to particular attributes. However, model specification often requires deep domain knowledge, as the functional form of the utilities needs to be specified a priori, often with simple and restrictive forms—such as linear in inputs and parameters to represent compensatory behavior. Furthermore, predictive performance of standard discrete choice tends to be lower than that of other approaches, such as ensemble models and deep neural networks [1].

Although not the primary focus of research in discrete choice, there has been increased interest in modeling social influence and peer effects on decision-making. Intuitively, if a large number of an individual’s friends purchase—or consider purchasing—a particular brand of cellphone, that individual may be inclined to do the same. In fact, decisions depend not only on personal preferences and the attributes of available options but also on how peers or social networks value those alternatives, a phenomenon known as network or peer effects. Nevertheless, discrete choice models often overlook network effects, and those that do account for them (e.g., [2], [3], [4], [5], [6], [7]) arguably lack the flexibility and predictive capabilities of other model classes, such as deep neural networks.

For instance, in [2], the author models binary transit choice in New York City using a linear specification for the latent utilities conditional on previous choices. This model depends on alternative attributes, socio-demographics, and mode shares in a previous time. In [3], the authors model career decisions after high school using a linear specification for the latent utilities that depend on alternative attributes, socio-demographics, and network latent utilities as an endogenous term. A more comprehensive model is proposed in [4], where the authors model the adoption of new technology using a linear specification for the latent utility associated with each alternative at every time step. This includes, apart from individual characteristics and alternative attributes, a term associated with previously acquired knowledge, a function of past individual decisions, and two types of social influence: (i) from individuals who previously chose the alternative and (ii) from those who did not.

In more theoretical approaches, such as the one described in [5], the authors introduce

network effects by including the strength of social utilities and a subjective expectation per individual about the share of people who make each choice. An important idea reflected in the model presented by the authors is that social interactions reflect what people think about the behavior of their peers (reference group), not necessarily their actual behavior. In their model, the authors impose that subjective beliefs are equal to conditional objective choice probabilities, which they assert are equivalent to rational expectations.

More recently, in [7], the authors examined the choice of the commute mode using a model that includes the spatial lag effect and spatially correlated error terms, similar to the SAL and SAE models from the spatial econometrics literature [8]. In addition, their model captures random taste variations with a spatially correlated structure. The most significant contribution of their research is that by allowing for taste variations to account for spatial correlations, the authors can account for residential self-selection. To our knowledge, their multinomial probit model is the most flexible and comprehensive specification that accounts for network effects in the current literature.

In all these models, social influence is incorporated into the latent utility as a weighted average of past or concurrent decisions, latent utilities, or expectations over reference group probabilities. This average is then multiplied by a parameter that reflects the overall significance of network effects. Furthermore, the specification of the latent utility is linear in terms of model parameters, alternative attributes, and socio-demographics. Therefore, to accommodate more complex forms of latent utility in these specifications, data pre-processing and feature engineering, expert knowledge, and behavioral assumptions are required, with the former requiring input from the latter two, which are not always available.

Deep learning (DL) models have recently been applied to discrete choice problems, often exhibiting higher predictive performance (e.g., [1], [9], [10]) and more flexibility, as universal approximators with automatic feature learning, than traditional DCMs. However, deep learning models are often regarded as black-box approaches with high on-sample predictive performance but limited potential for inference and, in some cases, poor generalization. Although this assertion may hold for general-purpose architectures, models that incorporate expert knowledge and inference frameworks accounting for parameter instability could potentially offer interpretability comparable to standard DCMs.

Recently, there has been growing interest in the deep learning literature regarding graph-based methods applied to tabular data. In [11], the authors present a brief survey of some of the available graph methods (e.g., Graph Neural Networks, Laplacian Regularization) applied to discrete choice problems. Graph Neural Networks have shown strong predictive performance, particularly in recent architectures like GCNII, as presented in [12], which leverage residual connections (as in ResNET [13]) and identity mappings. These architectures achieve the highest predictive performance on some benchmark datasets, second only to label graph-

based post-processing methods like the correct and smooth approach described in [14].

Recent research has demonstrated that deep learning models can be used to extract economic information, such as marginal substitution rates, resulting in insights that are as comprehensive as those obtained from standard models [15]. Furthermore, it has been shown that model ensembles, models with inductive biases, and training procedures incorporating gradient regularization [16] can produce insights consistent with behavioral intuition. However, despite these advancements, current interpretable deep learning models have yet to incorporate network effects.

In this study, we design an interpretable Graph Neural Network (GNN) model for discrete choice analysis that incorporates network effects and can induce independence from irrelevant alternatives (IIA) in the multinomial setting. To test our model, we use a dataset with revealed mode preferences, drawing from the 2010/2011 Regional Household Travel Survey conducted by the New York Metropolitan Transportation Council [17]. This dataset was extended with travel time and cost estimates obtained via the Google API, as well as social connection graphs created using trip origins and destinations from the Census Bureau [18], mapped to the census tracts recorded in the travel survey. We also use data from the 2016 U.S. election aggregated by county, along with a social connectedness graph, as in the case study presented in [11], to evaluate the performance of our model on problems with high class imbalance.

We compare our model’s predictive performance with that of traditional discrete choice models (DCMs), such as the standard logit model, and general-purpose deep learning architectures. Our findings demonstrate that our proposed architecture not only delivers behavioral insights that better align with intuition than those from off-the-shelf deep learning models, but also surpasses traditional DCMs in predictive performance. Furthermore, we implement approximate Bayesian inference using stochastic gradient Langevin dynamics (SGLD) [19] and demonstrate its potential to enhance the interpretability of general-purpose architectures within the context of our empirical example.

In the next section, 2, we provide a brief overview of standard DCMs, and then, in Section 3, we discuss the generalities of the current DCMs that incorporate spatial or network effects. In Section 4, we outline some guiding principles for designing interpretable deep learning models for discrete choice. Following this, in Section 5, we describe the GNN architecture that our model builds upon and that has gained popularity in the machine learning community. Then, in Section 6, we introduce our interpretable GNN that accounts for social influence in discrete choice problems. Subsequently, we discuss our case studies in Sections 7 and 8, and present the results from our model in Section 9. Finally, in Section 10, we offer concluding remarks and discuss new research avenues that we consider highly promising within the Deep Learning for Discrete Choice subfield.

2. Standard Discrete Choice Models

In the discrete choice framework employed in this work, we assume that when faced with the decision of choosing between two alternatives (e.g., public transit or a private car for commuting to work), an individual i will base their decision y_i on the maximization of their own utility u_i . This utility comprises a deterministic part that considers the individual’s characteristics and the attributes of the alternatives (differences in the binary setting), as well as a random part that captures unobservables.

In a standard binary discrete choice setting, the latent utility u_i that an individual i derives from choosing an alternative of interest over the second one (e.g., transit over car in a binary mode choice context) is represented by equation 1:

$$u_i = v_i + \epsilon_i, \tag{1}$$

where ϵ_i represents the random component used to account for unobserved effects, and v_i is the deterministic part of the utility, as given by equation 2:

$$v_i = \mathbf{x}_i^T \boldsymbol{\beta} + \mathbf{q}_i^T \boldsymbol{\gamma}. \tag{2}$$

Here, $\boldsymbol{\beta}$ and $\boldsymbol{\gamma}$ are estimable vectors of parameters, \mathbf{x}_i is a vector that represents the differences between the mode attributes (e.g., the difference in travel time), and \mathbf{q}_i is a vector of individual characteristics (e.g., household income). In basic binary choice models, the random component of the utility ϵ_i is typically assumed to be independent and identically distributed (i.i.d.) across decision makers, following either a standard normal distribution (binary probit model) or a standard logistic distribution (logistic regression or binary logit model¹).

In this binary setting, the probability of choosing the alternative of interest (i.e., $y_i = 1$, where y_i is a binary choice indicator) is given by equation 3:

$$P(y_i = 1) = P(\epsilon_i \leq v_i), \tag{3}$$

and the likelihood function for this model specification with n individuals is given by equation 4:

$$\mathcal{L}(\boldsymbol{\beta}, \boldsymbol{\gamma}) = \prod_{i=1}^n P(\epsilon_i \leq v_i)^{y_i} (1 - P(\epsilon_i \leq v_i))^{1-y_i}. \tag{4}$$

¹A binary logit model specifies a utility function for each alternative, with i.i.d. EV1 error terms. The likelihood of the binary logit model depends on the difference of the utilities, leading to u_i as specified above, and the difference of the two error terms, which, following the properties of the EV1 distribution, is logistically distributed, making the model in differences equivalent to a logistic regression.

For a more detailed and general description of standard discrete choice models, refer to [20].

3. Network and Peer Effects in Discrete Choice

The presence of network or peer effects in a discrete choice problem invalidates the i.i.d. assumption for the error terms in a standard discrete choice model. This violation leads to several issues, such as biased standard errors and, in cases where some covariates share the same correlation structure as the latent utility, even biased parameter estimates.

Network and peer effects have been extensively studied in spatial econometrics and, to some extent, in discrete choice analysis. In this section, we provide a brief review of the models in the literature that we consider to be the most widely adopted and foundational. However, this is not intended to be an exhaustive review.

3.1. Autoregressive Models from Spatial Econometrics

There are two general models that account for correlated observations applicable to both continuous and limited dependent variables [21], namely the Spatially Autoregressive Error and the Spatially Autoregressive Lag models. In this paper, we describe these approaches in the context of limited dependent variables, as encountered in discrete choice problems. The primary distinction between the two models lies in the assumption about the source of correlation in the continuous underlying index function (i.e., individual utilities in a discrete choice context). The Spatially Autoregressive Error (SAE) model assumes that correlation originates solely from correlated errors, whereas the Spatially Autoregressive Lag (SAL) model attributes autocorrelation to interactions between individual utilities. Both models result in a non-spherical variance-covariance structure capable of representing heteroskedasticity.

In the binary setting, the SAE model, which involves n individuals, k alternative attributes, and r socio-demographic variables, can be represented using the following matrix form for the latent utility, as shown in Equation 5:

$$\mathbf{u} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + (\mathbf{I} - \rho\mathbf{W})^{-1}\boldsymbol{\epsilon}, \quad (5)$$

where \mathbf{X} is an $n \times k$ matrix that contains alternative attribute differences for each individual, \mathbf{Q} is an $n \times r$ matrix that represents individual characteristics, \mathbf{I} is an $n \times n$ identity matrix, ρ is a scalar parameter that identifies the strength of correlated unobserved effects, \mathbf{W} is an $n \times n$ adjacency or connectivity matrix that represents the structure of interactions or connections between individuals or entities in a network, and $\boldsymbol{\epsilon}$ is an error vector that contains n uncorrelated elements. The adjacency matrix in peer effects econometrics is a foundational tool that encodes the structure of social or peer interactions, enabling researchers to quantify and analyze how individuals' behaviors or outcomes are shaped by their peers. In spatial

econometrics, where the focus is on geographic or spatial relationships between observations (e.g., regions or households), the matrix is commonly called the spatial adjacency matrix or simply the spatial weight matrix.

The SAL model, on the other hand, has the latent utility matrix specification shown in Equation 6:

$$\mathbf{u} = (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{Q}\boldsymbol{\gamma} + (\mathbf{I} - \rho\mathbf{W})^{-1}\boldsymbol{\epsilon}. \quad (6)$$

Social influence or peer effects are usually modeled using a framework akin to that of the Spatially Autoregressive Lag (SAL) model. To illustrate this, consider the latent utility function for the SAL model written in the following form:

$$\mathbf{u} = \rho\mathbf{W}\mathbf{u} + \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + \boldsymbol{\epsilon}. \quad (7)$$

In this formulation, $\rho\mathbf{W}\mathbf{u}$ represents the part of the latent utility influenced by peers, and ρ models the importance of peer effects. The term $\mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma}$ encapsulates the deterministic part of the private latent utility, which accounts for individual preferences and characteristics. Lastly, $\boldsymbol{\epsilon}$ denotes the uncorrelated error term, as in the SAE model.

3.2. State-of-the-Art DCMs with Network Effects

As previously discussed, the model proposed by [7] could be considered the most comprehensive specification for discrete choice models accounting for correlations between observations, which is a general Spatial Autoregressive with Autoregressive Disturbances (SARAR) model. This model works with panel data, models peer effects/social/network influence, and can account for self-selection effects by allowing for correlated random unobserved preference heterogeneity. A simplified binary specification of SARAR in matrix form, reminiscent of a more general form of the SAE and SAL models, follows the equation presented below:

$$\mathbf{u} = \rho\mathbf{W}\mathbf{u} + \mathbf{X}\left(b + (\mathbf{I} - \rho_{\beta}\mathbf{W})^{-1}\boldsymbol{\tau}_{\beta}\right) + \mathbf{Q}\boldsymbol{\gamma} + (\mathbf{I} - \rho_{\epsilon}\mathbf{W})^{-1}\boldsymbol{\epsilon} \quad (8)$$

In this specification, the social part of the utility is modeled as in the SAL model, and the correlated unobserved effects are modeled as in the SAE model. However, there are a couple of terms not present in either the SAE or SAL models. These terms are those used to account for self-selection effects, namely:

$$\boldsymbol{\beta} = b + (\mathbf{I} - \rho_{\beta}\mathbf{W})^{-1}\boldsymbol{\tau}_{\beta} \quad (9)$$

where $\boldsymbol{\tau}_{\beta}$ is an uncorrelated random vector of size n , b models the expected marginal utilities with respect to the alternative attributes, and ρ_{β} models the strength of self-selection based on attributes. With $\boldsymbol{\tau}_{\beta}$ and $\boldsymbol{\epsilon}$ as normally distributed random variables, this latent utility specification follows a normal distribution with a non-spherical variance-covariance matrix, as in the SAE and SAL models.

It is evident that restricting this SARAR specification allows us to recover either the SAE or SAL models, as well as the standard discrete choice model presented earlier. To illustrate this fact, setting $\rho = 0$ and $\rho_\beta = 0$ would yield the SAE model, and further setting $\rho_\epsilon = 0$ would yield the standard choice model discussed at the beginning of this paper.

This SARAR specification is fully transparent and interpretable; however, similar to standard discrete choice models, its specification limits the latent utilities to a linear function relative to alternative attributes and socio-demographics. Incorporating any additional complexity into the model would require expert knowledge as well as feature selection and engineering. Alternatively, we propose designing an interpretable graph neural network model. Such a strategy will enable us to retain the economic structure and information provided by the extant approaches while offering automatic feature learning and leveraging the predictive performance capabilities of deep learning.

3.3. Brief Note on Convergent Matrices and Affine Systems

In the specifications presented earlier, the term $(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{z}$ (where \mathbf{z} is a vector) appears multiple times, warranting a more in-depth discussion. Suppose that the individual latent utilities—or other individual-specific variables—are continuously updated according to the neighbors’ values and a constant term before the individuals decide which alternative to select. In discrete time, following the previously defined notation, this assumption can be represented as:

$$\mathbf{u}(t + 1) = \rho\mathbf{W}\mathbf{u}(t) + \mathbf{z} \quad (10)$$

Assuming that $\rho\mathbf{W}$ is a convergent matrix (which can be controlled by design), the sole equilibrium point for this system is given by:

$$\bar{\mathbf{u}} = (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{z} \quad (11)$$

Additionally, under the same assumption regarding $\rho\mathbf{W}$, the system reaches equilibrium in the limit:

$$\lim_{t \rightarrow \infty} \mathbf{u}(t) = \bar{\mathbf{u}} = (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{z} \quad (12)$$

Therefore, $(\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{z}$ represents the fixed point reached in the limit of the system defined by $\mathbf{u}(t + 1) = \rho\mathbf{W}\mathbf{u}(t) + \mathbf{z}$. For instance, with $\mathbf{z} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + \boldsymbol{\epsilon}$, the affine system

$$\mathbf{u}(t + 1) = \rho\mathbf{W}\mathbf{u}(t) + \mathbf{z}$$

reaches equilibrium at:

$$\mathbf{u} = (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{Q}\boldsymbol{\gamma} + (\mathbf{I} - \rho\mathbf{W})^{-1}\boldsymbol{\epsilon}, \quad (13)$$

which corresponds to the latent utility specification for the SAL model.

Although these facts might not seem immediately relevant, the reader will later find that they significantly clarify the application of Graph Neural Networks (GNNs) for discrete choice with social or network/peer effects.

4. Deep Learning for Discrete Choice

The primary limitation of deep learning models applied to discrete choice is their inherent lack of interpretability right out of the box and their tendency to overfit the training set, leading to overall poor generalization. In discrete choice research, inference is often considered more important than prediction. Therefore, applying general-purpose architectures from deep learning, such as fully connected neural networks, to discrete choice problems without meaningful architectural adaptations can result in economic insights that are not useful—a fact that has discouraged the research community.

The lack of interpretability in deep learning models applied to discrete choice stems from three main challenges, namely: (i) model architectures that lack behavioral interpretation, (ii) unstable parameter estimates, and (iii) a lack of methods for representing epistemic uncertainty. The first challenge was highlighted by [15], who demonstrated how to extract economic information as comprehensive as that provided by standard discrete choice models and illustrated methods for incorporating behavioral assumptions into model specifications through architectural design. The second challenge relates to the shape of the cost function for training when the model is underspecified by the data, as discussed by [22]. Specifically, when the model architecture involves a large number of parameters but is constrained by limited data, the cost function tends to exhibit extensive valleys and multiple local minima. This cost function shape can result in various parameter estimates with equivalent predictive performance on the training data, but with multiple economic implications. The third challenge-related closely to the second one—is associated with the fact that, in most inference problems, researchers work with limited data, leading to uncertainty in the hypotheses derived from the model-fitting process [23]. To our knowledge, the applications of deep learning models to discrete choice have not addressed the representation of epistemic uncertainty, so deep learning applications to discrete choice have focused solely on point estimation.

In this paper, we primarily focus on addressing the first challenge by designing an interpretable deep learning architecture to model discrete choice problems with network or social effects, exploiting convolutional graph neural networks. Additionally, we explore methods to address unstable parameter estimates and represent epistemic uncertainty by effectively implementing approximate Bayesian inference. However, the second and third challenges will be addressed in full in future work.

4.1. Stochastic Langevin Dynamics and Weight Averaging for Approximate Bayesian Inference

Recent applications of deep learning (DL) to discrete choice modeling have demonstrated that the interpretability of these models can be enhanced through the use of model ensembles [15] for inference, rather than relying solely on a single set of model weights \mathbf{w} . The effectiveness of this relatively straightforward technique stems from the fact that, for many problems, deep learning architectures tend to be underspecified by the training data. This underspecification leads to irregular loss landscapes characterized by large connected valleys and multiple modes [24, 25]. Therefore, for better interpretability, researchers are leaning towards approximate Bayesian approaches to deep learning, such as Stochastic Weight Averaging (SWA), Stochastic Weight Averaging Gaussian (SWAG) [24], and Stochastic Gradient Langevin Dynamics (SGLD) [19].

Stochastic Weight Averaging (SWA) is a computationally efficient way to deal with unstable parameter estimates. SWA offers better generalization and more stable solutions than those obtained from standard training methods while potentially providing better behavioral insights. The SWA approach averages the weight iterates from stochastic gradient descent during the learning process. The algorithm introduces minimal computational and memory overhead since the weight average is computed as a running average every c gradient updates. A simplified version of the learning procedure, presented in [24], is shown below:

Algorithm 1 Stochastic Weight Averaging [24]

Require: initial weights $\tilde{\mathbf{w}}$, cycle length c , number of epochs e , learning rate α , loss function L

- 1: $\mathbf{w} \leftarrow \tilde{\mathbf{w}}$ ▷ Initialize weights with $\tilde{\mathbf{w}}$
- 2: $\mathbf{w}_{SWA} \leftarrow \tilde{\mathbf{w}}$ ▷ Initialize SWA weights
- 3: $\eta_{models} \leftarrow 0$ ▷ Initialize number of models in average
- 4: **for** $i \leftarrow 1, 2, \dots, e$ **do**
- 5: $\mathbf{w} \leftarrow \mathbf{w} - \alpha \nabla L(\mathbf{w})$ ▷ Stochastic gradient update
- 6: **if** $\text{mod}(i, c) = 0$ **then**
- 7: $\mathbf{w}_{SWA} \leftarrow \frac{\mathbf{w}_{SWA} \cdot \eta_{models} + \mathbf{w}}{\eta_{models} + 1}$ ▷ Compute model average
- 8: $\eta_{models} \leftarrow \eta_{models} + 1$ ▷ Increment the number of models in the average

SWAG has a very similar implementation to SWA but differs in that it also computes the running weights' variance and uses these statistics to model weights as random variables drawn from a high-dimensional multivariate normal distribution. In that sense, SWA only provides a set of weights in a flatter region of the parameter space that has the potential for better generalization [24] but remains a single point estimate, while SWAG gives a true approximation to the weights' posterior that can be used for Bayesian inference.

Similar to SWA and SWAG, SGLD uses iterates from stochastic gradient descent. It approximates Langevin dynamics to obtain an approximation of the weights’ posterior distribution that can be treated as MCMC iterates and used for Bayesian inference. The algorithm works by injecting noise $\boldsymbol{\eta}_t$ into the (batched) stochastic gradient updates, as described by equations 14 and 15, and saving the weight iterates \boldsymbol{w}_t . In those equations, q_{ti} , x_{ti} , and y_{ti} represent the socio-demographics, alternative attributes, and selected alternative for individual i in batch t . $p(\boldsymbol{w}_t)$ is the prior distribution for the weights (e.g., ℓ_1 or ℓ_2 regularization in deep learning frameworks), and $p(q_{ti}, x_{ti}, y_{ti} | \boldsymbol{w}_t)$ is the likelihood of the observed individual under the set of weights \boldsymbol{w}_t . In Equation 14, N denotes the total number of training observations, n denotes the number of observations in the batch, and α_t denotes the gradient step size at t .

$$\Delta \boldsymbol{w}_t = \frac{\alpha_t}{2} \left(\nabla \log p(\boldsymbol{w}_t) + \frac{N}{n} \sum_{i=1}^n \nabla \log p(q_{ti}, x_{ti}, y_{ti} | \boldsymbol{w}_t) \right) + \boldsymbol{\eta}_t \quad (14)$$

$$\boldsymbol{\eta}_t \sim \mathcal{N}(0, \alpha_t) \quad (15)$$

With this gradient update structure, \boldsymbol{w}_t approaches samples from the posterior $p(\boldsymbol{w} | q_{ti}, x_{ti}, y_{ti})$, as these updates approximate Langevin dynamics, which converge to the posterior distribution [19]. In contrast with SWAG, SGLD does not approximate the posterior mode using a Gaussian distribution and does not require additional post-estimation sampling.

In this paper, we implement SGLD to perform approximate Bayesian inference using our models. However, as pointed out earlier, an exhaustive analysis of Bayesian deep learning for discrete choice, including interval estimation and hypothesis testing, is beyond the scope of this paper. We will address this in future work focusing on epistemic uncertainty representation.

5. One-size-fits-all Graph Convolutional Neural Networks

In this study, we use a type of Graph Neural Network (GNN) known as Graph Convolutional Neural Networks (GCNs), as detailed in Kipf and Welling (2016) [26]. GNNs are neural networks specifically designed to learn from graph-structured data, which includes tabular datasets with an underlying network structure. In the context of discrete choice, these network structures can be associated with social or geographic ties, as previously discussed. GCNs have been applied in various domains, including the analysis of physical systems, the prediction of protein interfaces, and the classification of diseases [27]. In discrete choice problems, GCNs have been used by [28] with standard architectures on benchmark datasets, with a focus on prediction capabilities.

To understand how GCNs are used within our discrete choice framework, consider the forward computation in a GCN layer, which is defined as:

$$\mathbf{A}^{(l+1)} = g^{(l+1)}\left(\mathbf{W} \mathbf{A}^{(l)} \Theta^{(l+1)}\right) \quad (16)$$

In this equation, $\mathbf{A}^{(l)} \in \mathbb{R}^{n \times o}$ denotes the hidden representation at layer l for all observations in the choice dataset, where n is the number of observations and o is the size chosen for the hidden representation at layer l . The activation function $g^{(l)}$ is a non-linearity (such as ReLU, as proposed in [26]), and $\Theta^{(l)} \in \mathbb{R}^{o \times p}$ is a matrix of learnable parameters that maps the hidden representations from \mathbb{R}^o to \mathbb{R}^p . As before, the matrix $\mathbf{W} \in \mathbb{R}^{n \times n}$ encodes the social connections or network structure.

If this operation is applied directly to the inputs of the neural network, the expression becomes:

$$\mathbf{A}^{(1)} = g^{(1)}\left(\mathbf{W} \text{CONCAT}(\mathbf{X}, \mathbf{Q}) \Theta^{(1)}\right). \quad (17)$$

Here, we substitute $\mathbf{A}^{(0)}$ with the operator $\text{CONCAT}(\mathbf{X}, \mathbf{Q})$, which concatenates socio-demographics and alternative attributes for each observation. When used in the output layer L , the computation then becomes:

$$\hat{\mathbf{y}} = \sigma\left(\mathbf{W} \mathbf{A}^{(L-1)} \Theta^{(L)}\right), \quad (18)$$

where $\hat{\mathbf{y}}$ is the vector of predictions and $\sigma(\cdot)$ is the logistic function (for the binary case). Since GCNs can be applied across any layer of the network, one could design a model consisting entirely of graph convolutional operations applied from the input to the output layers. For instance, a two-layer GCN model using the ReLU activation function would be given by:

$$\hat{\mathbf{y}} = \sigma\left(\mathbf{W} \text{ReLU}\left(\mathbf{W} \text{CONCAT}(\mathbf{X}, \mathbf{Q}) \Theta^{(1)}\right) \Theta^{(2)}\right). \quad (19)$$

This entire function is almost everywhere differentiable, allowing the model to provide economic insights comparable to those derived from standard discrete choice models. However, it is important to note that these insights may not always be behaviorally reasonable (e.g., not respecting monotonicity). Note that the network structure captured by \mathbf{W} is considered for each individual's choice prediction, and the model is no longer linear in inputs or parameters.

This GCN architecture is generic and could be applied to any binary classification problem with an underlying network structure. In this setup, network effects are intertwined with the private component of the deterministic utility (i.e., $\mathbf{X}\beta + \mathbf{Q}\gamma$ in a linear utility model), and the architecture does not directly enforce specific behavioral assumptions. In the following section, we build upon this foundation by designing a model specifically tailored for discrete choice, incorporating architectural design choices that are supported by the Machine Learning literature and ensure better behavioral interpretability.

6. GCNs Tailored for Discrete Choice: Skip-GNN

The model presented in the previous section represented the latent utility as:

$$\mathbf{u} = \mathbf{W} \mathbf{A}^{(L-1)} \Theta^{(L)}, \quad (20)$$

where $\mathbf{A}^{(L-1)} = f(\mathbf{X}, \mathbf{Q})$ is a non-linear function of socio-demographics and alternative attributes. We will depart from this general formulation of the latent utilities for our proposed architecture, integrating meaningful design choices based on two inductive biases that have proven empirically useful for estimation in the machine learning literature, namely: skip connections and batch normalization.

On the one hand, skip connections² [13] have been shown to be necessary for training deep learning models that are under-specified by the data. It was demonstrated in Li et al. (2018) [25] that architectures with skip connections have a dramatic effect on the loss landscape. The authors provide visualizations on random directions of the parameter space for different architectures and illustrate how skip connections prevent the loss landscape from exhibiting problematic levels of non-convexity. Their study elucidates the positive impact of skip connections on training speed and generalization, as observed empirically.

On the other hand, batch normalization, developed by Ioffe and Szegedy [29], was introduced to account for covariate shift (when the input distribution of a learning system changes) by normalizing the inputs to certain layers of deep networks trained with random batches of data. The inclusion of batch normalization in deep learning architectures has been shown to speed up the training process while having state-of-the-art performance –or better if combined with model ensembles [29].

In our architecture, we will use skip connections to distinguish between the linear and non-linear parts of the latent utilities, as well as between private and social utilities (the latter influenced by peers). Batch normalization will be implemented in one of the last layers of the model as a means to control the general scale and location of the latent utilities. In addition to incorporating behavioral insights, these inductive biases have the potential to produce model architectures that are easier to train. After training, these models are also likely to converge to generalize better. We name our architecture Skip-GNN because it relies heavily on skip and residual connections.

²In this paper, the term ‘skip connections’ is used to encompass both residual and skip connections. In both cases, information from previous layers is propagated forward in the network and is either concatenated or summed with the inputs of subsequent layers. These terms are often used interchangeably in the literature, their advantages during training and in terms of generalization are equivalent, and the distinction between them is not particularly relevant for the purposes of this paper.

6.1. Linear and Non-Linear Private Utility Components

The first assumption we make for our model, without any loss of generality, is the existence of an underlying private component in the latent utility that includes both linear and non-linear parts, as in a semi-parametric specification:

$$\mathbf{u}_{pr} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + f(\mathbf{X}, \mathbf{Q}), \quad (21)$$

where $f(\cdot)$ denotes a non-linear function, such as a fully connected neural network. This specific structure of the latent utilities can be efficiently implemented by employing skip connections from the input layer to the final layer of a deep learning architecture.

6.2. Private and Social Utilities

Next, to incorporate network effects in the model, consider the following system:

$$\begin{aligned} \mathbf{a}^{(l+1)} &= \mathbf{W}\mathbf{a}^{(l)}\theta^{(l)} + \mathbf{u}_{pr} \\ &= \mathbf{W}\mathbf{a}^{(l)}\theta^{(l)} + \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + f(\mathbf{X}, \mathbf{Q}). \end{aligned} \quad (22)$$

This system could be modeled using a Graph Convolutional Neural Network with one-dimensional latent representations $\mathbf{a}^{(l)} \in \mathbb{R}^{n \times 1}$, learnable parameters $\theta^{(l)} \in \mathbb{R}$, linear activation functions, and skip connections from the private utilities to each graph convolutional layer.

With a large number of layers L and setting $\theta^{(l)} = \rho, \forall l$ (such that $\mathbf{W}\mathbf{a}^{(l)}\theta^{(l)} = \rho\mathbf{W}\mathbf{a}^{(l)}$), and following the discussion on convergent matrices and affine systems from section 3, we can recover the following version of the SAL model:

$$\mathbf{u} = \mathbf{a}^L = (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{X}\boldsymbol{\beta} + (\mathbf{I} - \rho\mathbf{W})^{-1}\mathbf{Q}\boldsymbol{\gamma} + (\mathbf{I} - \rho\mathbf{W})^{-1}f(\mathbf{X}, \mathbf{Q}). \quad (23)$$

This specification would then incorporate a non-linear term that depends on alternative attributes and socio-demographics in the latent utility representation instead of the random term $\boldsymbol{\epsilon}$ from the original SAL model. Again, this structure for the latent utility can be conveniently implemented using skip connections.

Allowing the GCN to have non-linear mappings $g^{(l)}(\cdot)$, and for high-dimensional hidden representations $\mathbf{A}^{(l-1)}$, we would have:

$$\mathbf{A}^{(l)} = g^{(l)}(\mathbf{W}\mathbf{A}^{(l-1)}\boldsymbol{\theta}^{(l)} + \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + f(\mathbf{X}, \mathbf{Q})), \quad (24)$$

leading to the following latent utility model:

$$\mathbf{u} = \mathbf{W}\mathbf{A}^{(L-1)}\boldsymbol{\theta}^{(L-1)} + \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + f(\mathbf{X}, \mathbf{Q}). \quad (25)$$

The operations in Graph Convolutional layers can be interpreted as a social process where individuals update their utilities (represented as a latent high-dimensional embedding) in

discrete time. These representations, which may be high-dimensional, are based on their neighbors' representations and the parameters $\boldsymbol{\theta}^{(l)}$. In contrast to the SAL model and the SARAR model proposed by Bhat (2015) [7], the significance of peer effects is not encapsulated in a single parameter ρ , but rather in the set of parameters $\boldsymbol{\theta}^{(l)}$ associated with the GCN.

Additionally, this model can capture a form of exogenous interaction effects in the latent utilities, akin to those in the spatial Durbin model (see LeSage and Pace, 2009 [30]), by setting:

$$\mathbf{A}^{(0)} = \text{CONCAT}(\mathbf{X}, \mathbf{Q}), \quad (26)$$

and,

$$\mathbf{A}^{(l)} = \text{CONCAT}(g^{(l)}(\mathbf{W}\mathbf{A}^{(l-1)}\boldsymbol{\theta}^{(l)} + \mathbf{u}_{pr}), \mathbf{X}, \mathbf{Q}), \quad (27)$$

so that the operation $\mathbf{W}\mathbf{A}^{(l-1)}\boldsymbol{\theta}^{(l)}$ has terms of the form $\mathbf{W}\mathbf{X}\boldsymbol{\theta}_X^{(l)}$ and $\mathbf{W}\mathbf{Q}\boldsymbol{\theta}_Q^{(l)}$ that capture a form of exogenous interaction effects, as conceptualized in the spatial Durbin model. The learnable parameters are $\boldsymbol{\theta}^{(l)} \in \mathbb{R}^{1+k+r}$ for layers $l > 1$, and $\boldsymbol{\theta}^{(l)} \in \mathbb{R}^{k+r}$ for $l = 1$ (with k equal to the number of alternative attributes and r equal to the number of socio-demographics).

6.3. Setting the General Scale of Utilities Through Batch Normalization

Finally, batch normalization is employed in the non-linear part of the private utility $f(\mathbf{X}, \mathbf{Q})$. This means that, during training, batch statistics are computed for $f(\mathbf{X}, \mathbf{Q})$ to normalize their value according to the algorithm 2 presented below³. At prediction time, the whole sample statistics are used instead.

Algorithm 2 BatchNorm: Batch Normalizing Transform, applied to activation z over a mini-batch. [29]

Input: Values of z over a mini-batch: $B = \{z_1, \dots, z_m\}$

Output: $\{\hat{z}_i = \text{BN}_{\gamma, \beta}(z_i)\}$

$$\begin{aligned} \mu_B &\leftarrow \frac{1}{m} \sum_{i=1}^m z_i && \triangleright \text{mini-batch mean} \\ \sigma_B^2 &\leftarrow \frac{1}{m} \sum_{i=1}^m (z_i - \mu_B)^2 && \triangleright \text{mini-batch variance} \\ \hat{z}_i &\leftarrow \frac{z_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} && \triangleright \text{normalize} \end{aligned}$$

This operation is analogous to setting the location and general scale for the utilities in the standard logit and probit models [20]. With these design decisions, the general GCN model for discrete choice implemented in this paper comprises three main blocks: i) the linear part of the private utilities, ii) a fully-connected neural network $f(\mathbf{X}, \mathbf{Q})$ with normalized outputs

³In practice, the output \hat{z}_i is re-scaled by two learned parameters. In our context, this re-scaling operation would not determine the location or scale of the utilities. For our purposes, we ignore that re-scaling operation so that the batch normalization layer does not have any learnable parameters.

representing the non-linear part associated with the private utilities, and iii) a series of GCN layers used to model network effects and exogenous interaction effects. The entire model structure is summarized as follows:

$$\mathbf{u}_{pr} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Q}\boldsymbol{\gamma} + \text{BatchNorm}(f(\mathbf{X}, \mathbf{Q})) \quad (28)$$

$$\mathbf{A}^{(0)} = \text{CONCAT}(\mathbf{X}, \mathbf{Q}) \quad (29)$$

$$\mathbf{A}^{(l)} = \text{CONCAT}(g^{(l)}(\mathbf{W}\mathbf{A}^{(l-1)}\boldsymbol{\theta}^{(l)} + \mathbf{u}_{pr}), \mathbf{X}, \mathbf{Q}) \quad (30)$$

$$\mathbf{u} = \mathbf{a}^L = \mathbf{W}\mathbf{A}^{(L-1)}\boldsymbol{\theta}^{(L)} + \mathbf{u}_{pr} \quad (31)$$

$$\hat{\mathbf{y}} = \mathbf{p} = \sigma(\mathbf{u}), \quad (32)$$

where $\boldsymbol{\theta}^l \in \mathbb{R}^{1+k+r}$ for all GCN layers $l > 1$ (with k equal to the number of alternative attributes and r equal to the number of socio-demographics).

The binary model architecture is illustrated in Figure 1. In this figure, the private utilities are computed by passing the alternative attributes and socio-demographics ($\text{CONCAT}(X, Q)$) through a fully connected neural network (NN) with normalized outputs (BatchNorm) and, in parallel, a linear layer (Linear), and summing their outputs. The GCN blocks, up to $L - 1$, contain four sequential operations: i) GCN_l represents the $\mathbf{W}\mathbf{A}^{(l-1)}\boldsymbol{\theta}^{(l)}$ operation, ii) the output from GCN_l is summed with the private utilities \mathbf{u}_{pr} , iii) the result is passed through a ReLU non-linearity operator, and iv) the output from the activation is concatenated (CONCAT) with X and Q . The last GCN operation (GCN_L) outputs the socially-informed part of the utility. To obtain the latent utilities \mathbf{u} , the socially informed utilities are summed with the private utilities. Finally, to generate probability predictions \mathbf{p} in the binary problem setting, the latent utilities are passed through a sigmoid activation function. The part of the model within the dashed frame will be referred to here as a Binary Skip-GNN block, which is convenient for illustrating our model in the multinomial setting.

6.4. Multinomial Extension, IIA, and Shared Alternative Parameters

Our architecture can be extended for scenarios involving multiple alternatives. However, in such multinomial cases, a decision must be made regarding whether to ensure independence from irrelevant alternatives (IIA), as in [15], or allowing the model to actually learn substitution patterns from the data. In this subsection, we discuss both options.

The more straightforward case involves the model that ignores the Independence of Irrelevant Alternatives (IIA). In such cases, similar to standard discrete choice models, it is

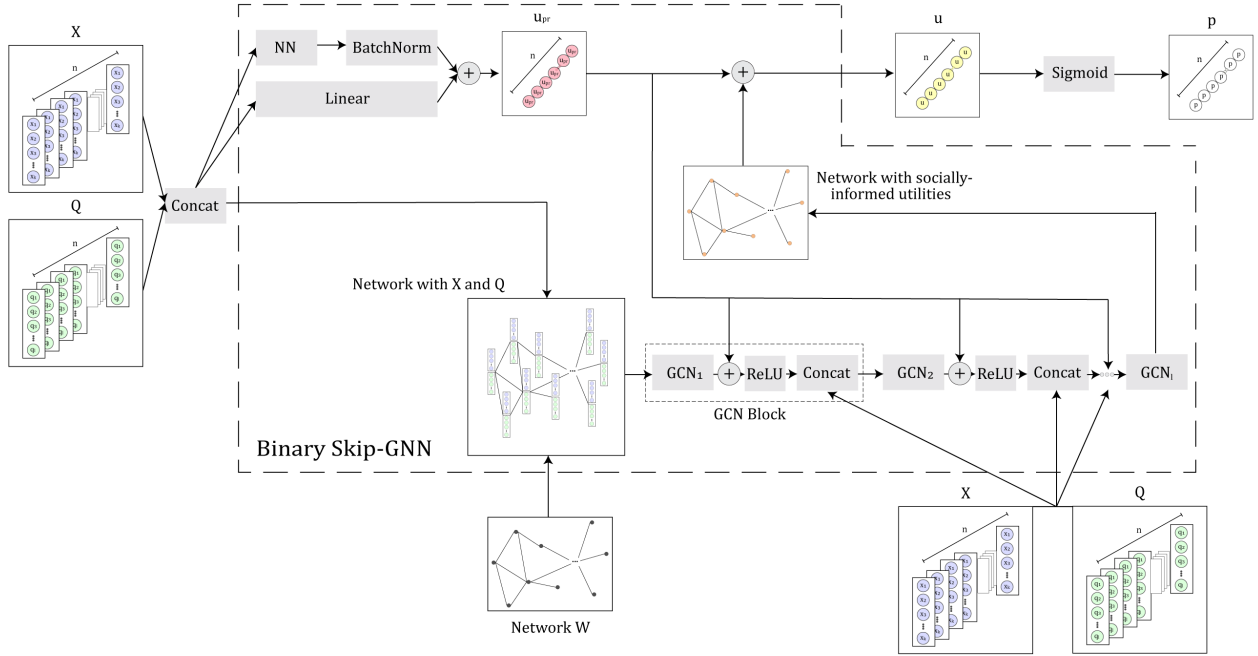


Figure 1: Binary Skip-GNN model architecture.

necessary to ensure that socio-demographic variables are incorporated into the utility functions for at most $J - 1$ alternatives (where J represents the total number of alternatives considered). If this is not ensured, the shape of the loss-function will include additional modes, affecting the stability of the parameters. Econometrically, the model would not be identified. Another important consideration pertains to non-linearity in the output layer. For the multinomial case, it is necessary to replace the Sigmoid function with a Softmax function to have an output size equal to the number of alternatives J , such that a full set of choice probabilities P_{ij} is produced. Other architectural adjustments are merely dimensional changes to accommodate J alternatives.

The more complex case would be ensuring the Independence of Irrelevant Alternatives (IIA). Whereas IIA is a reasonable axiom of preferences in certain theoretical contexts, it becomes problematic in logit-type models because IIA imposes restrictive and often unrealistic substitution patterns. However, IIA provides a baseline for checking regularity and ensures that predictions do not deviate too far from established economic theory, making it a useful tool for model validation and interpretability. One way to impose IIA employs 1D convolutions that effectively create independent alternative model blocks with shared parameters, as demonstrated in [15] and [31]. With these types of architectures, the original socio-demographics (common to all alternatives) cannot directly influence the alternative-specific utilities from the first layer, since the associated parameters are actually shared and

thus rendered irrelevant. To address this issue, the socio-demographics should go through another model (e.g., a fully connected neural network) before entering the alternative-specific utilities after a predetermined tunable number of layers. This model is used to create a socio-demographic embedding of size $J \times K$, so that the utilities for each alternative depend on K socio-demographic related features. This approach mirrors the one proposed in [31], with the key distinction being that the socio-demographic embeddings are not restricted to entering the utility function at the final layer. This additional flexibility allows for models that incorporate interactions between socio-demographics and alternative attributes.

The multinomial model architecture with IIA is depicted in Figure 2. As shown in that figure, the socio-demographics Q pass through a fully connected neural network (NN) that outputs J new socio-demographic vector representations Q_1, \dots, Q_J , one per alternative. The model has J binary skip-GNN blocks with shared parameters. The utilities u output by the J binary skip-GNNs go through a Softmax function to compute probability predictions p for each alternative.

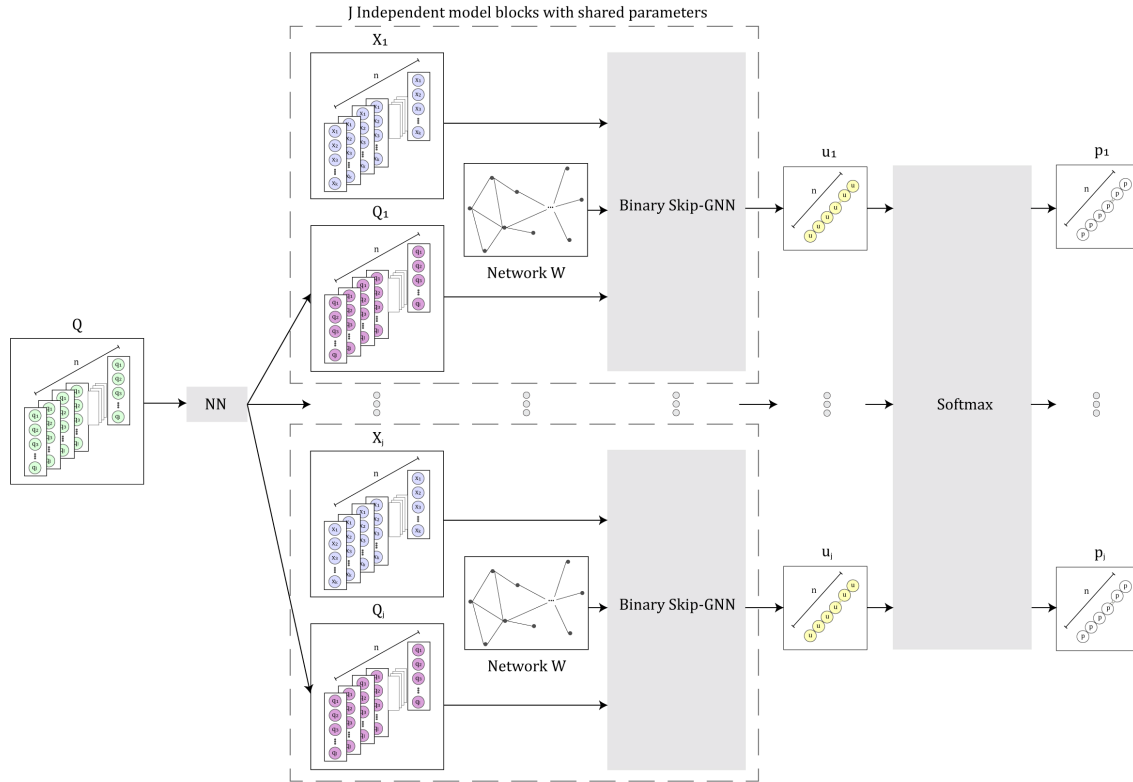


Figure 2: Multinomial Skip-GNN with IIA model architecture.

7. Mode Choice in New York City

We use the 2010/2011 Regional Household Travel Survey by the New York Metropolitan Transportation Council [17] to build mode choice data in New York City. The survey contains travel data from 18,965 households from areas in New York City, Long Island, the Hudson Valley, New Jersey, and Connecticut. The data is stored in a relational database with information on travel and mobility patterns, including socio-demographics, trip origins and destinations, and mode choices.

7.1. Binary Mode Choice: Public Transit vs Private Car

Following the reasoning presented in [2], we limit our analysis to the New York City areas. Additionally, we selected the home-based work (HBW) trips completed in a private car or transit longer than 1.5 miles, and we discarded observations with missing data. After the data selection and cleaning process, 2,444 trips are left for our analysis.

We use the US census block data [18] to obtain the longitude and latitude of origins and destinations for every trip. We used the Google API to retrieve estimates for each trip’s cost and travel time for various modes of transportation. This effort results in a novel binary mode choice dataset for NYC, containing trip origin and destination coordinates, trip costs, travel times, and socio-demographic variables (as presented in [32]). For our models, we consider the following variables: i) trip cost difference (transit vs. car), ii) travel time difference (transit vs. car), iii) an indicator for access to a private car in the household, iv) an indicator for the destination being in Manhattan, v) an indicator for high-income level, and vi) an indicator for declared gender. The description of these variables along with their mean values are presented in Table 1.

Figures 3 and 4 show the origin and destination locations (i.e. households and work locations), respectively. The red points represent the trips for which the individual selected a private car, and the blue points represent the ones for which the individual chose public transit. As shown in these plots, it is clear that there might be unobserved spatial effects that influence mode choice in New York City. For instance, for trips ending in Manhattan, individuals seem more likely to choose public transit over private cars.

In figures 5 and 6, we show trip examples for which public transit was selected, and in figures 7 and 8, we show trip examples for which private car was selected. The paths in blue correspond to the public transit routes, and the paths in red are for car. We show these figures to illustrate that mode choice is influenced by the alternative route characteristics, particularly the cost and time difference between modes.

Table 1: Summary of the variables considered for the binary mode choice problem.

Variable	Description	Mean
Trip cost difference	Transit cost minus private car cost (USD).	-3.36
Trip time difference	Transit travel time minus private car travel time (Minutes).	35.89
Vehicle availability	Indicator variable for car availability in the household. It takes the value of 1 if no cars are available in the household.	0.34
High income	Indicator variable for high income level ($> 100k$ USD per year).	0.33
Manhattan	Indicator variable for destinations in Manhattan.	0.47
Gender	Indicator variable for the male gender.	0.48
Mode choice (y)	Whether transit was selected as the travel mode. A choice indicator that takes the value of 1 if transit was selected.	0.61

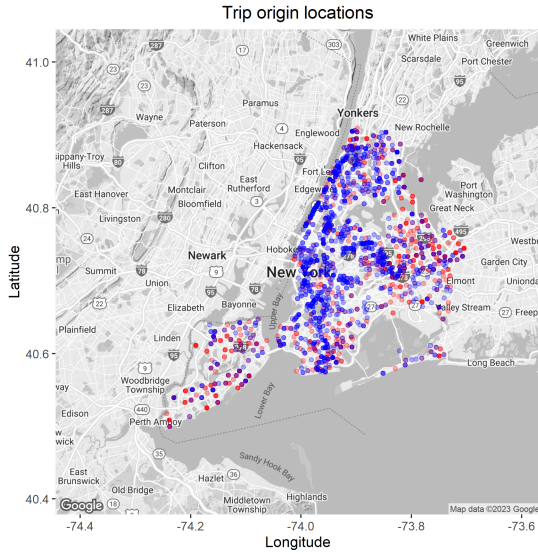


Figure 3: Trip origins and mode choice in New York City.

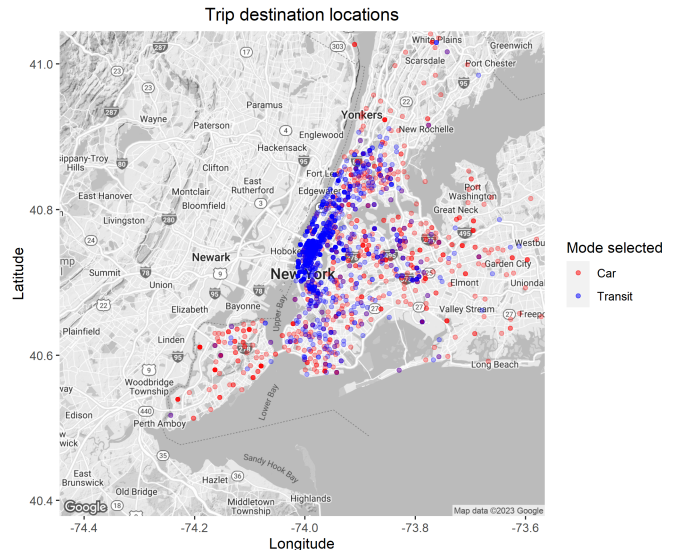


Figure 4: Trip destinations and mode choice in New York City.

7.2. Multinomial Mode Choice: Public Transit, Private Car and Non-motorized

For the multinomial case, we also limit our analysis to the New York City areas. We select home-based work (HBW) trips completed in transit, private car, or non-motorized modes (i.e. walking and bicycling) without limiting the trip length. Trip cost for the non-motorized alternative is set to zero and the travel time is computed using the Google API

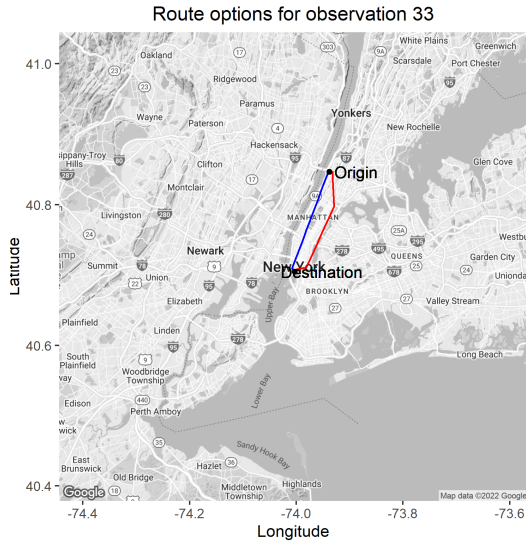


Figure 5: Trip 33 alternative routes for public transit and private car. For this trip, the cost of using a private car is \$4.75 higher than the cost of using public transportation.

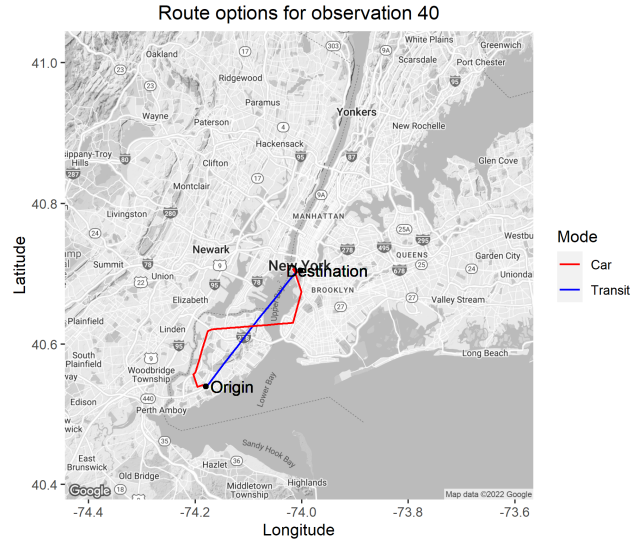


Figure 6: Trip 40 alternative routes for public transit and private car. For this trip, the cost of using a private car is \$11.77 higher than the cost of using public transportation.

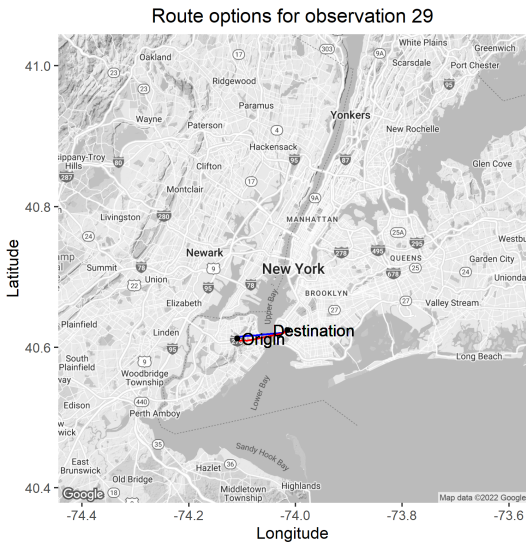


Figure 7: Trip 29 alternative routes for public transit and private car. For this trip, the cost of using a private car is \$1.06 higher than the cost of using public transportation.

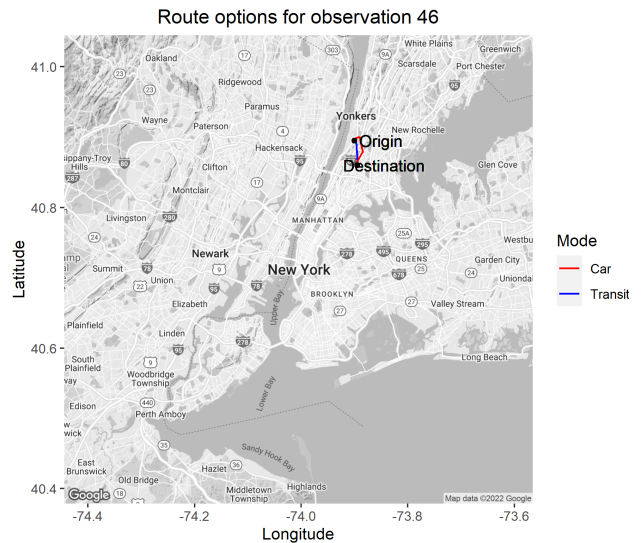


Figure 8: Trip 46 alternative routes for public transit and private car. For this trip, the cost of using a private car is \$0.10 higher than the cost of using public transportation.

for the cases where it is not revealed. We discarded observations with missing data, so after the data selection and cleaning process, 3,277 trips are left for our analysis. The description

for the variables considered in the multinomial problem is presented in Table 2. As for the binary problem, this dataset was also used in [32].

Table 2: Summary of the variables considered for the multinomial mode choice problem.

Variable	Description	Mean
Trip cost transit	Transit cost (USD).	2.35
Trip cost car	Car cost (USD).	4.64
Trip cost non-motorized	Non-motorized cost (USD).	0.00
Trip time transit	Transit travel time (Minutes).	50.38
Trip time car	Car travel time (Minutes).	18.32
Trip time non-motorized	Non-motorized travel time (Minutes).	93.85
Vehicle availability	Indicator variable for car availability in the household. It takes the value of 1 if no cars are available in the household.	0.36
High income	Indicator variable for high income level ($> 100k$ USD per year).	0.32
Manhattan	Indicator variable for destinations in Manhattan.	0.46
Gender	Indicator variable for the male gender.	0.47
Car mode share	Proportion of trips completed by car	0.37
Transit mode share	Proportion of trips completed by transit	0.51
Non-motorized mode share	Proportion of trips completed by non-motorized means of transportation	0.12

8. US county election 2016

As second case study, we use the 2016 county election dataset as presented in [11]. The data includes socio-demographic characteristics (e.g., death and birth rates, net migration, median income) aggregated by county, election results, and a network constructed using the Social Connectedness Index, which measures the relative frequency of Facebook friendships between each pair of counties. We use this dataset in a binary discrete choice setting with two candidate options, namely: Hillary Clinton and Donald Trump. All other candidates are excluded for simplicity, as no county selected a candidate outside of these two. In total, there are 3112 counties included in the dataset.

The map of the election results by county is presented in Figure 9. As shown, nearly 85% of counties had a majority of the population voting for Trump. This figure makes evident that neighboring counties often voted similarly. This is particularly noticeable in the blue clusters on the map, where several adjacent counties supported the Democrat candidate.

Table 3: Summary of the variables considered for the binary election choice problem.

Variable	Description	Mean
Death Rate	Percentage of deaths in the population annually.	10.81
Birth Rate	Percentage of births in the population annually.	11.62
Net Migration Rate	Percentage change in population due to migration (inflow minus outflow).	-0.04
Bachelor Rate	Percentage of adults aged 25 and older with at least a bachelor’s degree.	21.56
Median Income	Median household income in 2016 US dollars.	49403
Unemployment Rate	Percentage of the civilian labor force that is unemployed and seeking employment.	5.20
Rural-urban Continuum Code 2013	Classification of counties based on population density and proximity to metropolitan areas (1 = most urban, 9 = most rural). Dummy encoded and first category dropped.	-
Economic Typology 2015	Classification of counties based on predominant economic activity (e.g., farming, manufacturing, mining). Dummy encoded and first category dropped.	-
Election Choice (y)	Binary variable indicating whether Donald Trump (1) or Hillary Clinton (0) received the majority vote in the county.	0.8425

9. Results and Discussion

In this section, we evaluate and compare the predictive performance and behavioral insights obtained from our proposed architecture against those derived from an out-of-the-shelf GNN and a traditional logit model. The comparison includes results from both the standard models and their Stochastic Gradient Langevin Dynamics (SGLD) counterparts.

For the deep learning models, we perform a random grid search for the layer width, the number of fully connected neural network layers, the number of GCN layers, the weight decay rate, and the learning rate. We use ReLU activations for all non-linearities except for the output layer, for which we use the Sigmoid activation (binary) or Softmax (multinomial). For SGLD, we run the algorithm for 100,000 epochs with thinning every 1,000 epochs to reduce MCMC iterate autocorrelation.

2016 U.S. Presidential Election Results by County

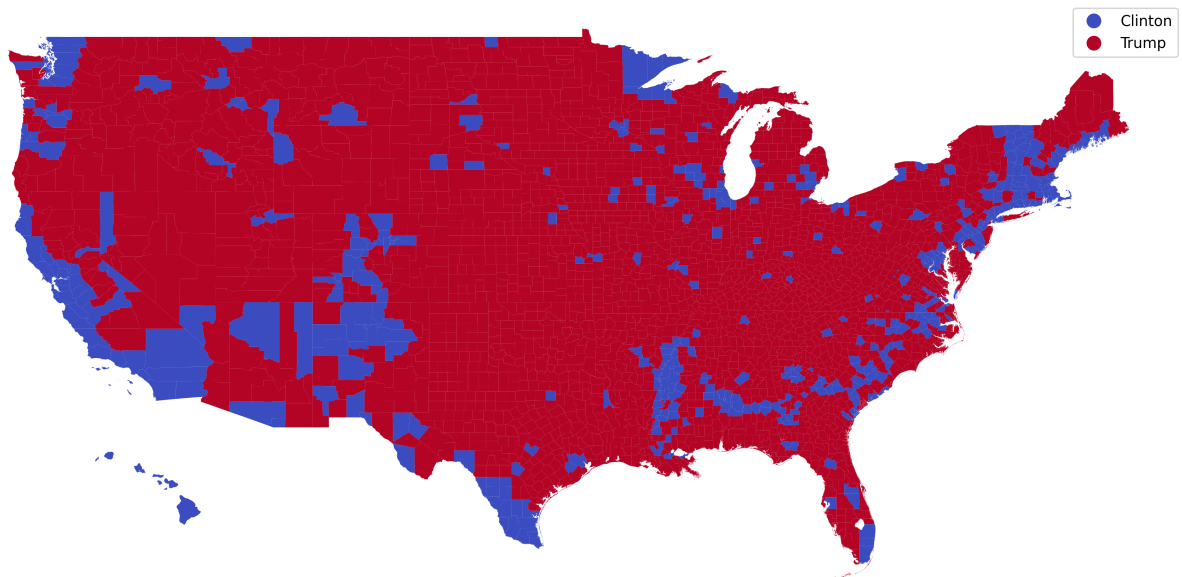


Figure 9: 2016 US Election Results by County.

9.1. Case study: Binary Mode Choice in NYC

In this subsection, we present the marginal utilities and value of travel time savings VOTTs estimated for the binary NYC mode choice dataset using a standard logit model, our Skip-GNN architecture, and a general-purpose GNN. We demonstrate that our Skip-GNN model provides insights that align with behavioral intuition, while SGLD (Stochastic Weight Averaging) improves the behavioral alignment of the off-the-shelf GNN.

9.1.1. Marginal Utilities

For the logit model, we found the marginal utility of travel time, denoted as β_t , to be -0.022, and the marginal utility of trip cost, denoted as β_c , to be -0.101. These negative values align with behavioral expectations, as anticipated. In the case of the general-purpose graph neural network, the individual marginal utility of travel time and trip cost are illustrated in Figures 10 and 12, respectively. For our proposed Skip-GNN model, the corresponding marginal utilities are detailed in Figures 11 and 13.

Regarding the marginal utilities of travel time, it is observed that the general-purpose graph neural network (GNN) yields positive marginal utilities for a significant proportion of individuals (as shown in Figure 10), which contradicts micro-economic expectation. Nonetheless, the median marginal utility of travel time across the sample aligns with the expected sign, being equal to -0.02. The discrepancies between the Stochastic Gradient Langevin

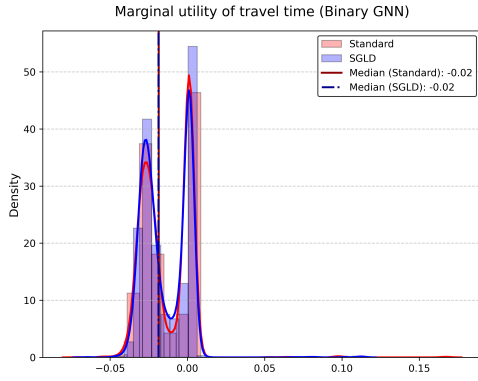


Figure 10: Marginal utility of travel time from general purpose GNN.

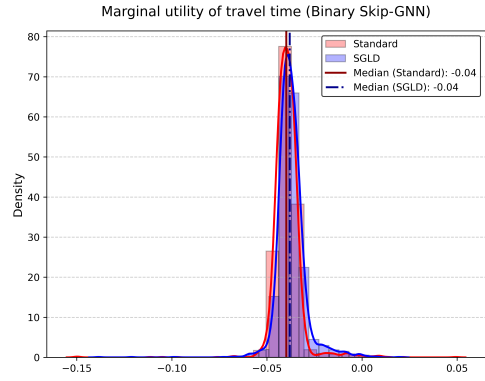


Figure 11: Marginal utility of travel time from the Skip-GNN.

Dynamics (SGLD) approach and the regular estimation procedure in this context are not pronounced. Conversely, for our proposed skip-GNN model, the marginal utilities of travel time are consistently negative across the entire sample (as shown in Figure 11), thereby enhancing regularity of the outcome.

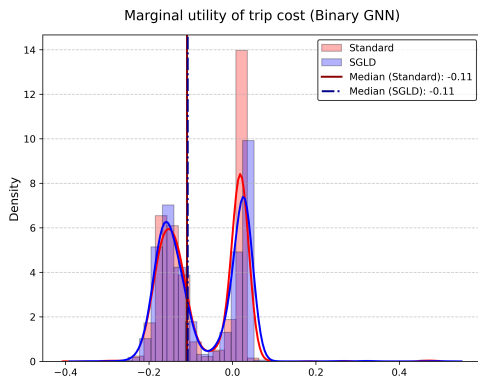


Figure 12: Marginal utility of trip cost for the general purpose GNN.

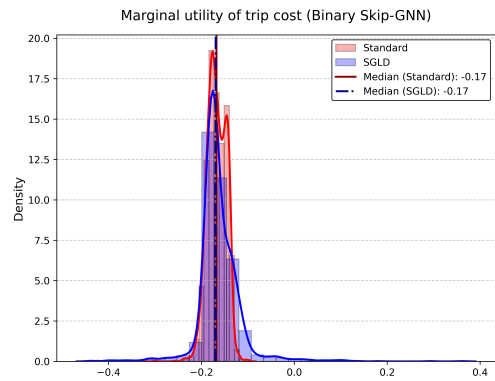


Figure 13: Marginal utility of trip cost for the Skip-GNN.

For the marginal utility of trip cost, we observe similar results as those just discussed for the marginal utility of travel time. The general-purpose architecture yields positive values for a substantial proportion of individuals (as indicated in Figure 12), in contrast to the consistently negative values generated by our proposed skip-GNN model across the dataset (Figure 13). Nevertheless, similar to the findings for travel time, the median marginal utility of trip cost derived from the general-purpose GNN aligns with the expected negative marginal

effect, being -0.11.

9.1.2. Value of Travel Time Savings

Regarding the value of travel time savings (VOTT), which is the marginal rate of substitution between travel time and travel cost and thus expected to be positive (when representing the willingness to pay to reduce travel time by a marginal unit), the point estimate obtained using the logit model is 15.75 USD per hour reduction in travel time. The VOTT values for the entire sample as derived from the general-purpose graph neural network (GNN) are illustrated in Figure 14, while those obtained from our skip-GNN model are presented in Figure 15. These figures provide insights into the distribution or average values of VOTT, as determined by the respective models.

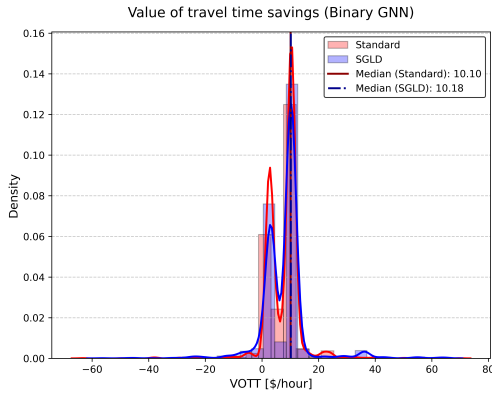


Figure 14: value of travel time savings from the general purpose GNN.

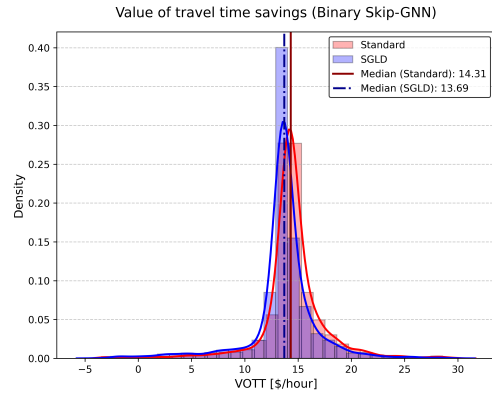


Figure 15: value of travel time savings from the Skip-GNN.

For the standard GNN, the median VOTT across the sample is calculated to be 10.10 USD per hour. In contrast, for the SGLD variant of this model, the VOTT is slightly higher, at 10.18 USD per hour. As depicted in Figure 14, the distribution of VOTT estimates for the standard GNN model includes a significant proportion of individuals with negative VOTT values, with values ranging from -70 USD to 80 USD. However, the SGLD version shows a more concentrated distribution around the mean with a reduced frequency of negative VOTT values.

The median VOTT estimated using our skip-GNN model is 14.31 USD per hour for the conventional estimates. In the case of the SGLD variant, the median VOTT is slightly lower, at 13.69 USD per hour. As illustrated in Figure 14, the VOTT estimates for the entire sample using both the regular and SGLD versions of the skip-GNN model exhibit very

sensible values.⁴

9.2. Case Study: Multinomial Mode Choice in NYC

In this subsection, we present the marginal utilities and values of time estimated for the multinomial mode choice dataset using a standard conditional logit model, our Skip-GNN architecture (both with and without representing the Independence of Irrelevant Alternatives (IIA)), and a general-purpose GNN. We demonstrate that our proposed model offers insights that align more closely with behavioral intuition than those provided by the general-purpose GNN. Additionally, we observe that SGLD have a notable effect on the marginal utilities and VOTT for our Skip-GNN model.

9.2.1. Marginal Utilities

We discuss first the histograms for the marginal utilities of travel time and trip cost for transit and private car for the whole sample of individuals. The marginal utilities for car travel time and trip cost, obtained from the general-purpose GNN model, are illustrated in Figures 16 and 17, respectively. The marginal utilities for transit travel time and trip cost are illustrated in Figures 18 and 19. As shown in these figures, the marginal utilities exhibit both positive and negative values across the sample at seemingly equal proportions. This observation contradicts micro-economic sign expectation, which would suggest that travel mode utilities should be lower for alternatives with higher costs or longer trip times.

The marginal utilities for car travel time and trip cost, derived from our Skip-GNN model, are depicted in Figures 20 and 21, respectively. Similarly, the marginal utilities for the transit mode are shown in Figures 22 and 23. Contrasting these results to those from the general-purpose architecture, our findings reveal marginal utilities that are more consistent with behavioral expectations, as evidenced by histograms that show very few observations with positive marginal utilities.

The marginal utilities for car travel time and trip cost, obtained from our Skip-GNN-IIA model, are illustrated in Figures 24 and 25, respectively. Similarly, the marginal utilities for the transit mode are displayed in Figures 26 and 27. The results from this variant, which applies the Independence of Irrelevant Alternatives (IIA) restriction, are on par with those from the Skip-GNN model without IIA. This consistency indicates that our model can operate effectively with or without the IIA constraint, without compromising its ability to align

⁴In [33], the authors found an average value of travel time savings (VOTT) for the low-income population to be 21.67 USD/hour, for the not low-income population to be 28.05 USD/hour, for the student population to be 10.96 USD/hour, and the average VOTT for the senior population to be 10.93 USD/hour. Their estimates are based on a group-level agent-based mixed (GLAM) logit applied to synthetic nation-wide level data. To the best of our knowledge, there are no other recent estimates derived from discrete choice models for the value of travel time savings in NYC.

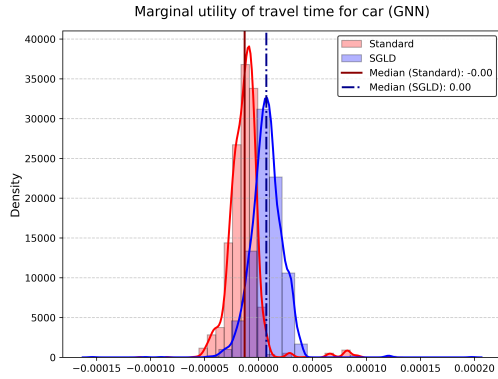


Figure 16: Marginal utility of travel time for car from general purpose GNN.

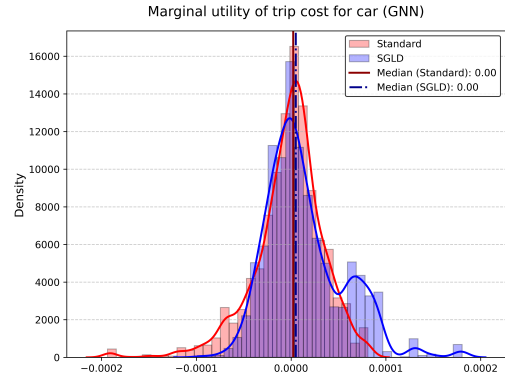


Figure 17: Marginal utility of trip cost for car from general purpose GNN.

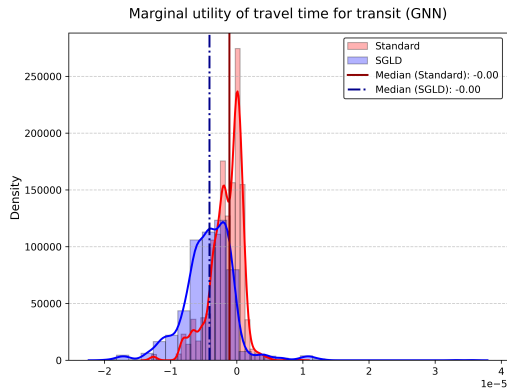


Figure 18: Marginal utility of travel time for transit from general purpose GNN.

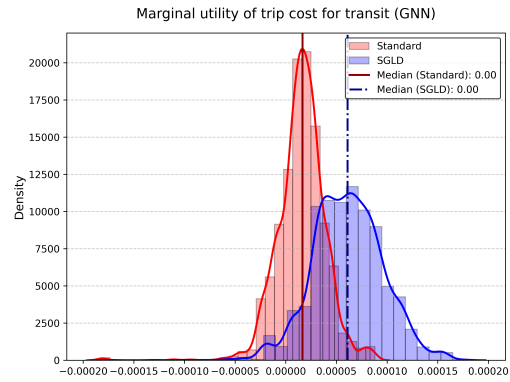


Figure 19: Marginal utility of trip cost for transit from general purpose GNN.

with expected behavioral regularity⁵. It is also important to note that the marginal utilities obtained with SGLD tend to be further away from zero than those found using the regular estimation procedure, which translates to a lower—or even null—proportion of individuals with marginal utilities that defy behavioral intuition.

9.2.2. Value of Travel Time Savings

The estimated value of travel time savings (VOTT), using a standard conditional logit model for the multinomial mode choice problem, was determined to be 12.64 USD. This estimate, along with the mean VOTT estimates from our models and the general-purpose GNN,

⁵For model selection, when behavioral regularity is met, out-of-sample prediction metrics should be used.

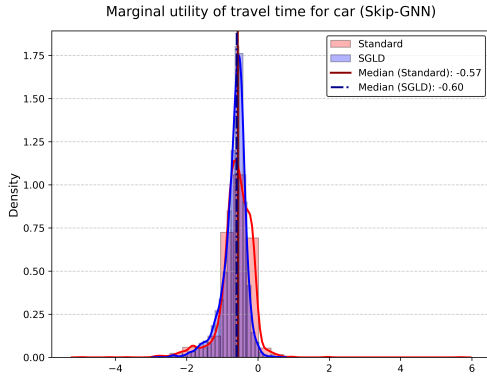


Figure 20: Marginal utility of travel time for car from Skip-GNN model.

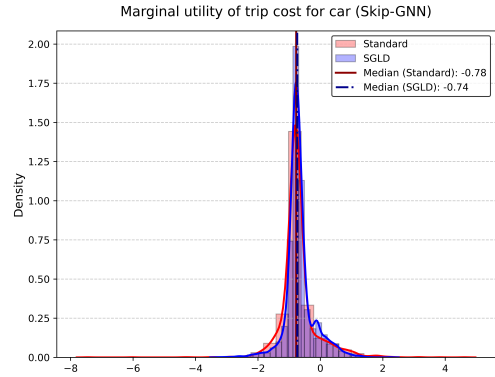


Figure 21: Marginal utility of trip cost for car from Skip-GNN model.

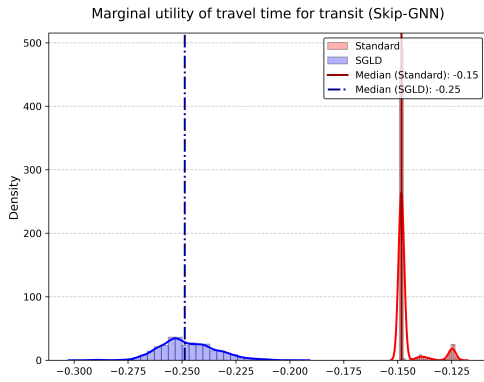


Figure 22: Marginal utility of travel time for transit from Skip-GNN model.

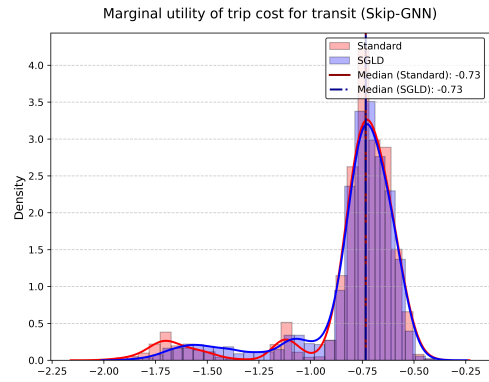


Figure 23: Marginal utility of trip cost for transit from Skip-GNN model.

are summarized in Table 4. Our proposed models yield reasonable median VOTT estimates. In contrast, the general-purpose deep learning architecture produced negative median VOTT estimates for both modes. These findings, along with the ones for the marginal utilities presented earlier, highlight the potential of our models to provide more plausible estimates than the ones that could be generated from off-the-shelf deep learning models.

In Figures 28 through 33, we present the histograms of the value of travel time savings (VOTT) for car and transit modes as estimated by the general-purpose GNN and our two model variants. For the general-purpose GNN, a significant number of individuals exhibit negative VOTTs, which contrasts with the consistently positive VOTT estimates for the majority of the sample produced by our models. This result aligns with our expectations, as we have previously discussed that our models' marginal utilities are consistent with behavioral

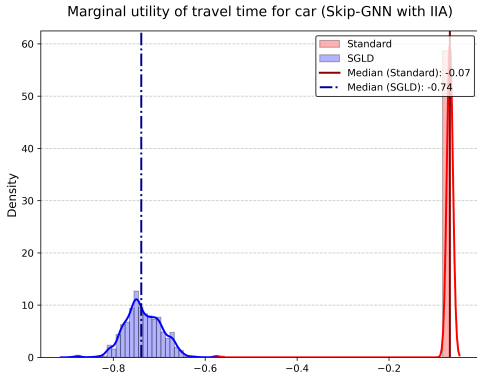


Figure 24: Marginal utility of travel time for car from Skip-GNN-IIA model.

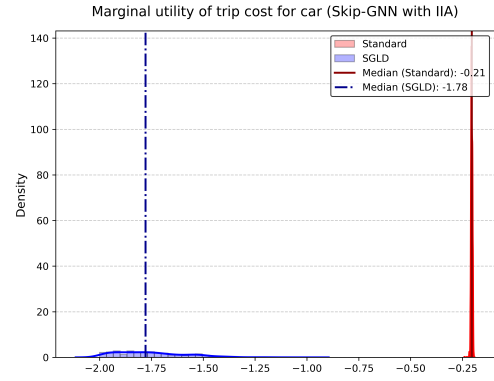


Figure 25: Marginal utility of trip cost for car from Skip-GNN-IIA model.

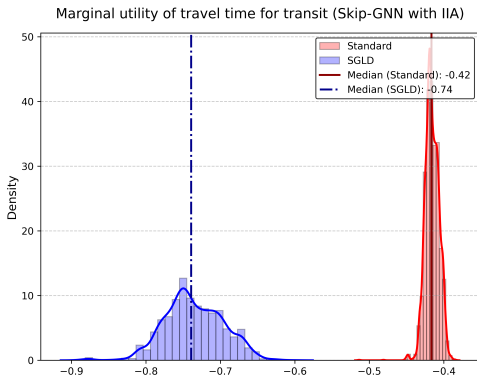


Figure 26: Marginal utility of travel time for transit from Skip-GNN-IIA model.

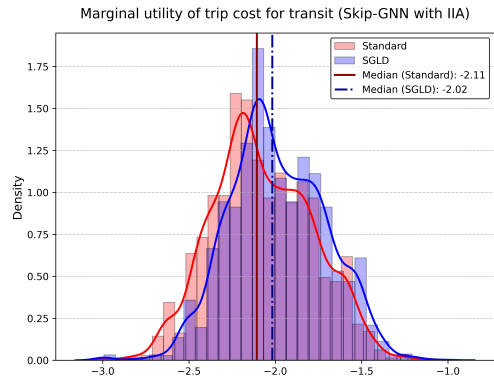


Figure 27: Marginal utility of trip cost for transit from Skip-GNN-IIA model.

intuition.

9.3. Case Study: Binary US Election

Inference for the binary U.S. election is conducted with respect to median income, unemployment rate, and bachelor's degree rate. We compute the odds ratios for these variables across all models and compare their values. For logit models, the odds ratios are constant and given by $\exp(\hat{\beta}_k)$, where $\hat{\beta}_k$ is the estimated parameter associated with the variable of interest. For the GNN and Skip-GNN models, odds ratios are computed from marginal effects using the partial derivative of the representative utilities with respect to the socio-demographic x_k . Since the odds ratios for the GNN and Skip-GNN models depend on the values of the socio-demographic vector \mathbf{x} , we compute the odds ratios for each individual, calculate their

Table 4: Median Value of travel time savings [USD/hour].

	Conditional Logit	GNN	Skip-GNN	Skip-GNN-IIA
VOTT Car	12.64	-9.03	35.54	19.72
VOTT Transit	12.64	-2.61	12.12	11.95

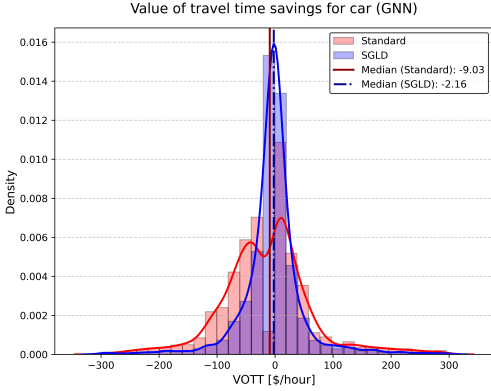


Figure 28: Value of travel time savings for car from the general purpose GNN.

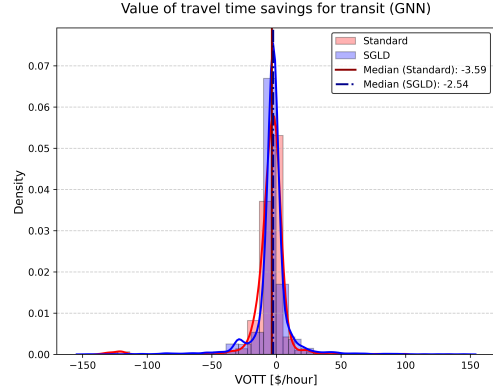


Figure 29: Value of travel time savings for transit from the general purpose GNN.

means and medians across the sample, and present their histograms.

9.3.1. Median Income Odds Ratio

Using the logit model, we found an odds ratio of 1.10 for a \$1000 USD increase in county median income. This implies that with a \$1000 USD increase, the odds of voting for Trump increase by 10% according to the logit model. In Table 5, we present the mean and median values of the odds ratio for median income found using the GNN and Skip-GNN models.

Table 5: Median Income Odds Ratio.

	GNN	GNN (SGLD)	Skip-GNN	Skip-GNN (SGLD)
Mean	1.14	1.28	1.13	1.15
Median	1.15	1.28	1.10	1.17

The odds ratio histogram for the GNN model is presented in Figure 34, while the histogram for the Skip-GNN model is presented in Figure 35.

For both models, the odds ratios are generally above 1.00, which indicates that an increase in the average county income increases the odds of voting for Trump in most counties.

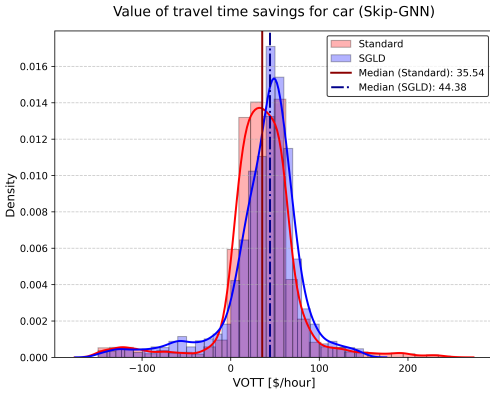


Figure 30: value of travel time savings for car from the Skip-GNN model.

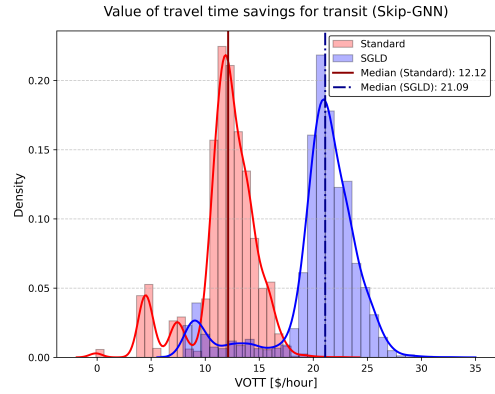


Figure 31: value of travel time savings for transit from the Skip-GNN model.

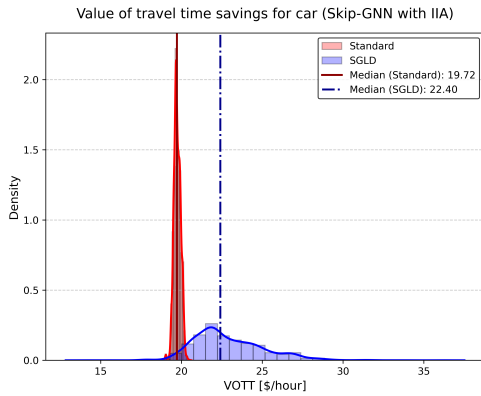


Figure 32: value of travel time savings for car from the Skip-GNN-IIA model.

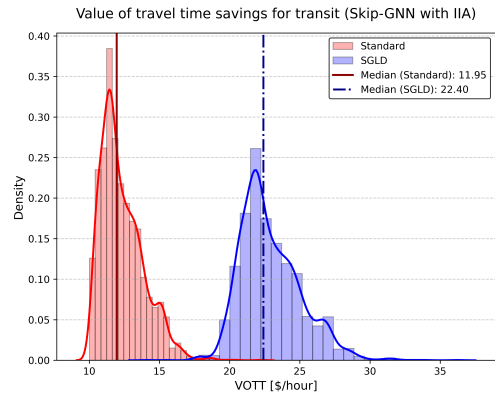


Figure 33: value of travel time savings for transit from the Skip-GNN-IIA model.

9.3.2. Bachelor Degree Rate Odds Ratio

Under the logit model, the odds ratio for a 1% increase in the bachelor's degree rate is 0.82. This means that a 1% increase in the percentage of people with a bachelor's degree in a county decreases the odds of voting for Trump by 18%. In Table 6, we present the mean and median values of the odds ratio for an increase of 1% in the bachelor rate found using the GNN and Skip-GNN models.

The odds ratio histogram for the GNN model is presented in Figure 36, while the histogram for the Skip-GNN model is presented in Figure 37.

For both models, the odds ratios are generally below 1.00, which indicates that an increase in the bachelor's degree rate reduces the odds of voting for Trump in most counties.

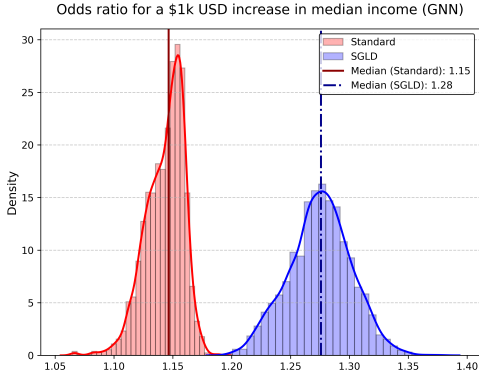


Figure 34: Odds ratio for an income increase of \$1000 USD under the GNN model.

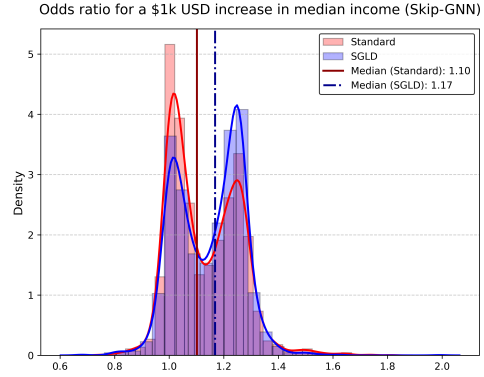


Figure 35: Odds ratio for an income increase of \$1000 USD under the Skip-GNN model.

Table 6: Bachelor Rate Odds Ratio.

	GNN	GNN (SGLD)	Skip-GNN	Skip-GNN (SGLD)
Mean	0.70	0.53	0.69	0.68
Median	0.70	0.54	0.71	0.66

9.3.3. Unemployment Rate Odds Ratio

The odds ratio for a 1% increase in the unemployment rate under the logit model is 0.64. Therefore, a 1% increase in the unemployment rate decreases the odds of voting for Trump by 36%. In Table 7, we present the mean and median values of the odds ratio for an increase of 1% in the unemployment rate found using the GNN and Skip-GNN models

Table 7: Unemployment Rate Odds Ratio.

	GNN	GNN (SGLD)	Skip-GNN	Skip-GNN (SGLD)
Mean	0.67	0.66	0.53	0.57
Median	0.66	0.66	0.53	0.54

The histogram for the GNN model is presented in Figure 38, while the one for the Skip-GNN model is presented in Figure 39.

As with the case, under both models, the results indicate that with a 1% increase in the unemployment rate, the odds of voting for Trump decrease in most counties.

9.4. Model Accuracy across Case Studies

To estimate the prediction performance out of sample of all models, we employed 5-fold cross-validation. In all cases, we compute the weighted (or balanced) accuracy, which is the

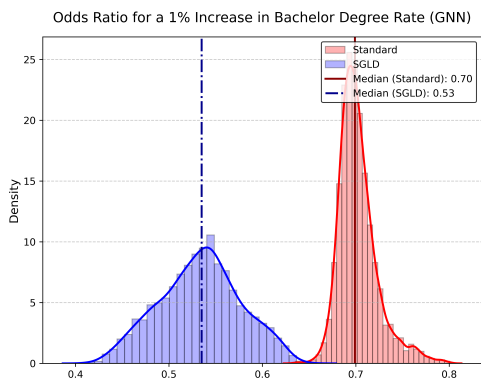


Figure 36: Odds ratio for an increase of 1% in the bachelor rate under the GNN model.

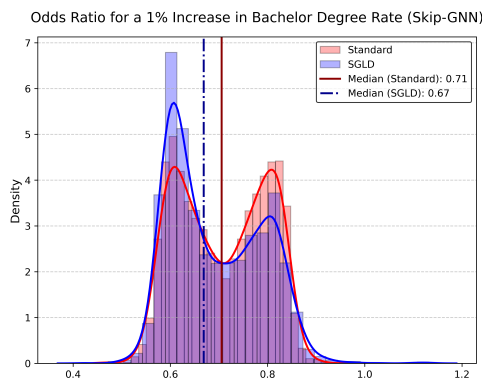


Figure 37: Odds ratio for for an increase of 1% in the bachelor rate under the Skip-GNN model.

average of specificity and sensitivity, also called the true positive and true negative rates for binary classification problems. In the multinomial mode choice problem, the weighted accuracy is calculated as the average per-class precision. This metric provides a more reliable measure of predictive performance as it accounts for class imbalances.

We observed a very modest improvement in prediction accuracy for our Skip-GNN model compared to the logit model in the binary mode choice scenario, and around a 6 percentage point increase in the multinomial mode choice setting (for both the IIA and unconstrained versions of our model). For the binary election dataset, we observed an improvement of more than 6 percentage points for our Skip-GNN model compared to the standard logit⁶. The general-purpose GNN is outperformed by our Skip-GNN model across all datasets, and for the binary mode choice problem, it does not even achieve a prediction performance on par with the logit model.

These gains align with the findings of [1], which reported an average performance improvement of around 5 percentage points for deep neural networks over traditional Discrete Choice Models (DCMs). For our case studies, it is important to note that even for the mode choice problems where the logit model outperforms the general-purpose GNN architecture, our Skip-GNN model is capable of attaining the highest out-of-sample performance. Detailed out-of-sample weighted accuracy metrics for all estimated models are presented in Table 8.

⁶To ensure a fair comparison against the deep learning models that are estimated using a weighted binary cross-entropy loss function, we estimated the logit models using per-class weights in the log-likelihood. The weighted accuracy computed for the US 2016 binary election with the logit model, without class weights, was estimated to be 72.68%. For the mode choice problems, the weighted accuracies with and without class weights are similar to one another.

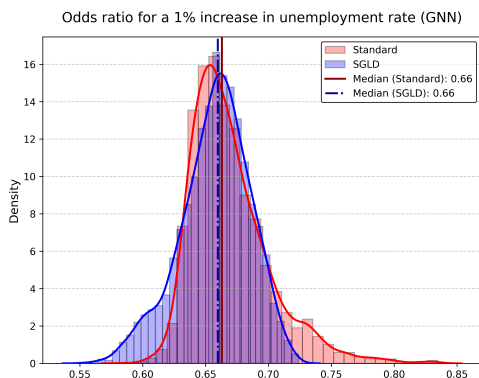


Figure 38: Odds ratio for an increase of 1% in the unemployment rate under the GNN model.

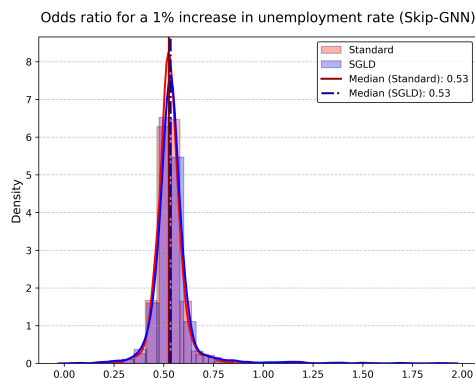


Figure 39: Odds ratio for for an increase of 1% in the unemployment rate under the Skip-GNN model.

Table 8: Accuracy on the test set model comparison

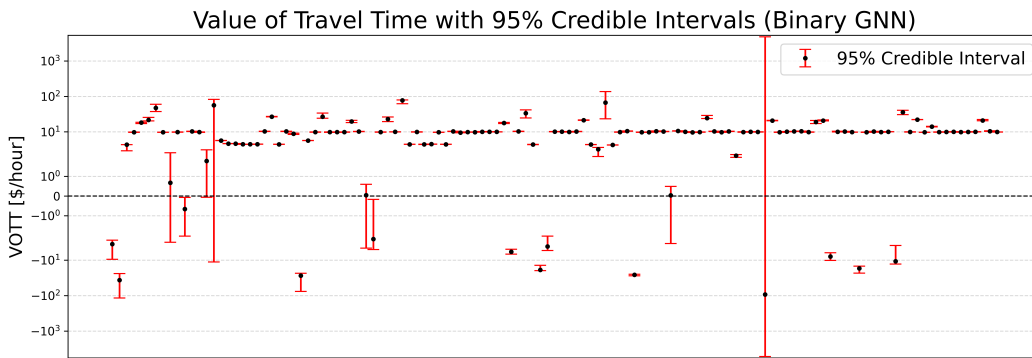
	Logit	Skip-GNN (IIA)	Skip-GNN	GNN
Binary mode choice	83.92%	NA	84.53%	76.75%
Multinomial mode choice	72.66%	77.70%	78.45%	64.36%
US 2016 binary election	81.21%	NA	87.29%	81.87%

9.5. Individual-Level Inference

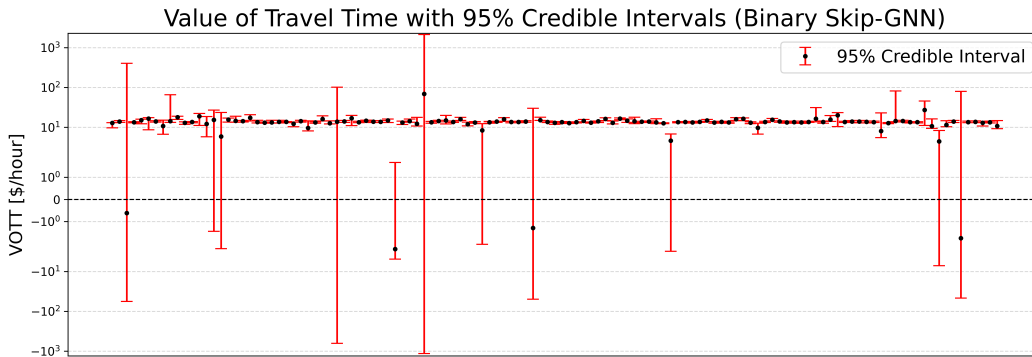
One of the advantages of deep learning models for discrete choice is the associate ability to perform inference at the individual level. For instance, in the binary case study presented in this paper, we estimate the Value of Travel Time (VOTT) for all individuals in the dataset, considering their socio-demographics and the attributes of the alternatives presented to them. Traditionally, in discrete choice modeling, this level of granularity is typically achieved using models that incorporate random preference heterogeneity, such as mixed-logit models, or by designing latent utility functions to capture deterministic preference heterogeneity based on socio-demographic interactions with alternative attributes. However, deep learning models achieve individual-level estimates through automatic feature learning, leveraging hidden interactions between input variables.

Using SGLD, we are thus able to provide individual-level estimates for the VOTT as well as 95% credible intervals. In Figure 40a, we present the credible intervals at the individual level for the VOTT for the GNN model for a sub-sample of 123 individuals, and in Figure 40b, we show the corresponding intervals for our Skip-GNN model. It is important to note

that for the GNN model, there are credible intervals with upper and lower limits that are negative. This indicates, with high posterior probability, that the VOTT for these individuals is negative, a result that defies behavioral intuition. In contrast, for our Skip-GNN model, while some individual-level VOTT median values are negative, the credible intervals for those individuals include both positive and negative values. This implies that there is insufficient posterior evidence under our model to suggest that the VOTT for those individuals is negative.



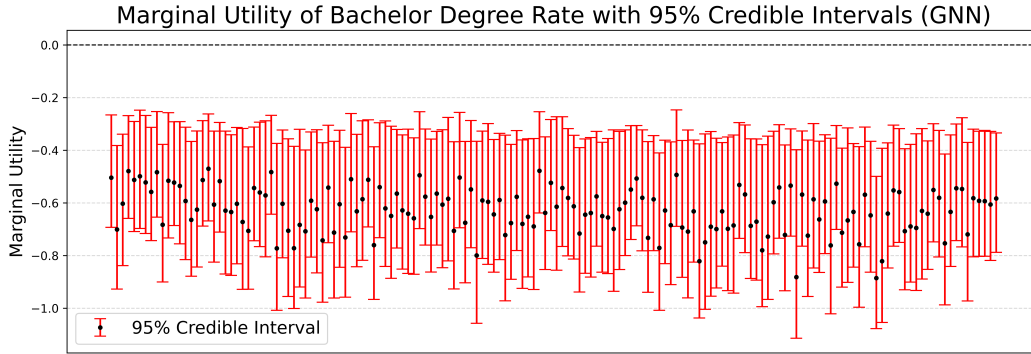
(a) Credible intervals for individual-level VOTT for binary GNN. Subsample of 123 observations.



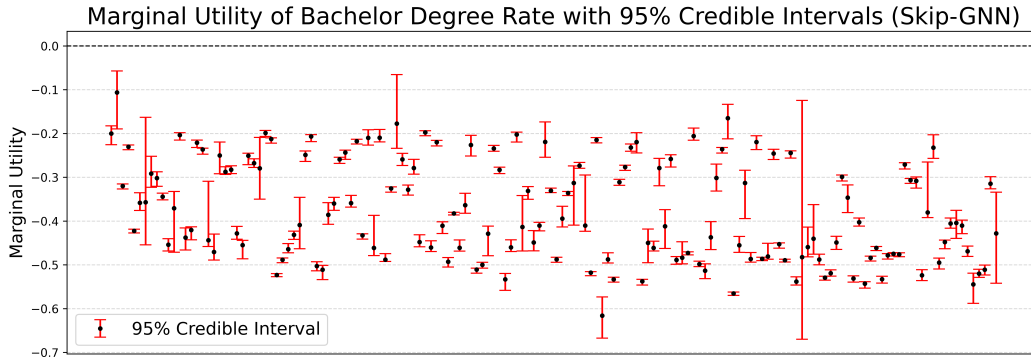
(b) Credible intervals for individual-level VOTT for binary Skip-GNN. Subsample of 123 observations.

Figure 40: Credible intervals for individual-level VOTT for the binary GNN and Skip-GNN models. Subsample of 123 observations.

As another example, in Figures 41a and 41b we provide the county-level credible intervals for the bachelor rate odds-ratios for both the GNN and Skip-GNN models, respectively. In this case, the credible intervals for the GNN model are evidently wider than the ones for our Skip-GNN model and the credible intervals overlap.



(a) Credible intervals for county-level bachelor rate odds-ratio for GNN. Subsample of 156 observations.



(b) Credible intervals for county-level bachelor rate odds-ratio for Skip-GNN. Subsample of 156 observations.

Figure 41: County-level credible intervals for bachelor rate odds-ratios using GNN and Skip-GNN models.

As mentioned earlier, this paper does not examine the representation of epistemic uncertainty for constructing credible intervals and performing hypothesis testing. Instead, the aim of this section is to highlight the importance of interval estimation and demonstrate how it enables researchers to draw meaningful conclusions and compare insights derived from competing models. Future work will provide a more comprehensive comparison of methods for representing epistemic uncertainty with Bayesian and not Bayesian methods, and evaluate their empirical coverage through simulation analysis.

10. Conclusions

In this paper, we have introduced a novel and interpretable GNN architecture for discrete choices, that we have named Skip-GNN, that captures social influence and exogenous interaction effects. Our model design enables interactions between attributes and socio-demographics through automatic feature learning. Our design decisions set the general scale

and location of utilities and clearly distinguish the individual’s private and socially informed components of utilities. Importantly, we observe that our approach enhances the behavioral intuition of the econometric parameter estimates without relying on hard or soft gradient constraints.

We have presented our model making sure to provide justifications for every design decision based on both discrete choice microeconomic theory and empirical findings from the machine learning literature. We offer more than just a useful deep learning architecture capable of modeling social influence and representing the Independence of Irrelevant Alternatives (IIA) if needed; we also equip discrete choice modelers with design principles. For example, the use of Batch Normalization (BatchNorm) helps set the scale and location of utilities. We have demonstrated that these principles can be used to develop models that not only deliver superior behavioral insights compared to off-the-shelf deep learning models but also attain high predictive performance.

We provide a binary version of our model and two alternatives for multinomial problems: one that reflects IIA and another without that constraint (which is a straightforward generalization of the binary model). For the IIA scenario, we build on ideas previously presented by Wang (2020) [15] and address socio-demographic variables by creating independent embeddings for each alternative, following a similar approach to that in Arkoudi (2023) [31]. We have called this IIA variant Skip-GNN-IIA. Both our Skip-GNN and Skip-GNN-IIA models exhibit high predictive performance and align with behavioral intuition, demonstrating that our model is suitable for either modeling flexible substitution patterns or for being restricted by IIA.

We tested our models on mode choice data from NYC in both binary and multinomial settings. Our models attain the highest out-of-sample prediction performance and exhibit strong alignment with behavioral intuition. For instance, in estimating the marginal utilities of travel time and trip cost, which are attributes that when increased make alternatives less attractive, our architecture consistently produces negative values across the sample of individuals. This contrasts with the outcomes from a general-purpose deep learning model which, despite presenting negative median values for the sample, reveals a substantial proportion of individuals with positive marginal utilities for travel time and trip cost. This observation, coupled with the detailed interpretation of each design decision, underscores the importance of tailoring architectures to provide more reliable and interpretable behavioral insights.

We also used 2016 binary election data from the U.S., aggregated by county, along with a network constructed using Facebook friendship data [11]. We observed a significant improvement in the weighted accuracy of our Skip-GNN model compared to both standard GNNs and traditional logit models. Leveraging the interpretability of our models, we computed estimates for the odds ratios associated with an increase in median income, bachelor’s de-

gree rate, and unemployment rate. Our findings suggest that an increase in median income raises the odds of voting for Trump, whereas an increase in the bachelor’s degree rate and unemployment rate decreases them in most counties.

While providing insights aligned with behavioral intuition, we also observe an increase in out-of-sample weighted accuracy for our models compared to standard DCMs and off-the-shelf GCNs across all three datasets, with a 6-percentage-point improvement in both the binary election data and the multinomial mode choice dataset. These results align with previous empirical studies showing performance gains when applying deep learning to discrete choice settings [1].

We used Stochastic Gradient Langevin Dynamics (SGLD) to account for epistemic uncertainty. SGLD enabled us to provide interval estimates for the value of travel time savings (VOTT) and odds ratios at the individual level. To our knowledge, this study is the first to provide interval estimates for a deep learning model applied to discrete choices. In future work, we will focus more comprehensively on interval estimation and hypothesis testing using approximate Bayesian inference.

Future work will thus focus on two main areas, namely: (i) using approximate Bayesian approaches in deep learning to better represent the epistemic uncertainty associated with discrete choice models, and (ii) exploring the connection between the private component of utility, as presented in this paper, and Gaussian Processes. Ultimately, we believe that incorporating deep learning into discrete choice modeling can be effectively achieved by carefully designing architectures (such as the one presented here), approaching the estimation problem from a Bayesian perspective, or—most likely—a combination of both.

Acknowledgments

This research was supported by the National Science Foundation Award No. SES-2342215. We are also thankful for the financial support provided by the Fulbright Scholarship Program, which is sponsored by the U.S. Department of State, the Colombian Fulbright Commission, and the Colombian Science Ministry. This project is solely the responsibility of the authors and does not necessarily represent the official views of the Fulbright Program, the U.S. government, or the Colombian government.

References

- [1] S. Wang, B. Mo, S. Hess, J. Zhao, Comparing hundreds of machine learning classifiers and discrete choice models in predicting travel behavior: an empirical benchmark, arXiv preprint arXiv:2102.01130 (2021).
- [2] F. Goetzke, Network effects in public transit use: evidence from a spatially autoregressive mode choice model for new york, *Urban Studies* 45 (2008) 407–417.
- [3] B. Li, R. Sickles, J. Williams, Estimating peer effects on career choice: A spatial multinomial logit approach, in: *Essays in Honor of Cheng Hsiao*, Emerald Publishing Limited, 2020.
- [4] A. Páez, D. M. Scott, E. Volz, A discrete-choice approach to modeling social influence on individual decision making, *Environment and Planning B: Planning and Design* 35 (2008) 1055–1069.
- [5] W. Brock, S. N. Durlauf, *Multinomial choice with social interactions*, 2003.
- [6] E. R. Dugundji, J. L. Walker, Discrete choice with social and spatial network interdependencies: an empirical example using mixed generalized extreme value models with field and panel effects, *Transportation Research Record* 1921 (2005) 70–78.
- [7] C. Bhat, A new spatial (social) interaction discrete choice model accommodating for unobserved effects due to endogenous network formation, *Transportation* 42 (2015) 879–914.
- [8] L. Anselin, *Spatial econometrics, A companion to theoretical econometrics* (2001).
- [9] Y. Han, F. C. Pereira, M. Ben-Akiva, C. Zegras, A neural-embedded discrete choice model: Learning taste representation with strengthened interpretability, *Transportation Research Part B: Methodological* 163 (2022) 166–186.
- [10] J. Lu, Y. Meng, H. Timmermans, A. Zhang, Modeling hesitancy in airport choice: A comparison of discrete choice and machine learning methods, *Transportation Research Part A: Policy and Practice* 147 (2021) 230–250.
- [11] K. Tomlinson, A. R. Benson, Graph-based methods for discrete choice, *Network Science* 12 (2024) 21–40. doi:10.1017/nws.2023.20.
- [12] M. Chen, Z. Wei, Z. Huang, B. Ding, Y. Li, Simple and deep graph convolutional networks, in: *International conference on machine learning*, PMLR, 2020, pp. 1725–1735.

- [13] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recognition, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 770–778.
- [14] Q. Huang, H. He, A. Singh, S.-N. Lim, A. R. Benson, Combining label propagation and simple models out-performs graph neural networks, arXiv preprint arXiv:2010.13993 (2020).
- [15] S. Wang, B. Mo, J. Zhao, Deep neural networks for choice analysis: Architecture design with alternative-specific utility functions, Transportation Research Part C: Emerging Technologies 112 (2020) 234–251.
- [16] S. Feng, R. Yao, S. Hess, R. A. Daziano, T. Brathwaite, J. Walker, S. Wang, Deep neural networks for choice analysis: Enhancing behavioral regularity with gradient regularization, arXiv preprint arXiv:2404.14701 (2024).
- [17] 2010/2011 Regional Household Travel Survey, Technical Report, New York Metropolitan Transportation Council, North Jersey Transportation Planning Authority, 2014.
- [18] Census tracts, <https://www.census.gov/geographies/reference-files/time-series/geo/gazetteer-files.2010.html#list-tab-264479560>, Accessed: 2023.
- [19] M. Welling, Y. W. Teh, Bayesian learning via stochastic gradient langevin dynamics, in: Proceedings of the 28th international conference on machine learning (ICML-11), Citeseer, 2011, pp. 681–688.
- [20] K. E. Train, Discrete Choice Methods with Simulation, 2 ed., Cambridge University Press, 2009. doi:10.1017/CB09780511805271.
- [21] L. Anselin, R. Florax, S. J. Rey, Advances in spatial econometrics: methodology, tools and applications, Springer Science & Business Media, 2013.
- [22] A. G. Wilson, P. Izmailov, Bayesian deep learning and a probabilistic perspective of generalization, Advances in neural information processing systems 33 (2020) 4697–4708.
- [23] A. G. Wilson, The case for Bayesian deep learning, arXiv preprint arXiv:2001.10995 (2020).
- [24] P. Izmailov, D. Podoprikin, T. Garipov, D. Vetrov, A. G. Wilson, Averaging weights leads to wider optima and better generalization, arXiv preprint arXiv:1803.05407 (2018).
- [25] H. Li, Z. Xu, G. Taylor, C. Studer, T. Goldstein, Visualizing the loss landscape of neural nets, Advances in neural information processing systems 31 (2018).

- [26] T. N. Kipf, M. Welling, Semi-supervised classification with graph convolutional networks, arXiv preprint arXiv:1609.02907 (2016).
- [27] J. Zhou, G. Cui, S. Hu, Z. Zhang, C. Yang, Z. Liu, L. Wang, C. Li, M. Sun, Graph neural networks: A review of methods and applications, *AI Open* 1 (2020) 57–81.
- [28] K. Tomlinson, A. R. Benson, Graph-based methods for discrete choice, arXiv preprint arXiv:2205.11365 (2022).
- [29] S. Ioffe, C. Szegedy, Batch normalization: Accelerating deep network training by reducing internal covariate shift, in: *International conference on machine learning*, pmlr, 2015, pp. 448–456.
- [30] J. LeSage, R. K. Pace, *Introduction to spatial econometrics*, Chapman and Hall/CRC, 2009.
- [31] I. Arkoudi, R. Krueger, C. L. Azevedo, F. C. Pereira, Combining discrete choice models and neural networks through embeddings: Formulation, interpretability and performance, *Transportation research part B: methodological* 175 (2023) 102783.
- [32] D. F. Villarraga, R. A. Daziano, Hierarchical nearest neighbor gaussian process models for discrete choice: Mode choice in new york city, *Transportation Research Part B: Methodological* 191 (2025) 103132. URL: <https://www.sciencedirect.com/science/article/pii/S019126152400256X>. doi:<https://doi.org/10.1016/j.trb.2024.103132>.
- [33] X. Ren, J. Y. Chow, Block-group level mode choice parameters for new york city and new york state [supporting dataset], the TRIS and ITRD database (2023).