

Investigating Links between Illicit Massage Businesses through Natural Language Processing and Graph Machine Learning

Anonymous ACL submission

Abstract

Human trafficking exploits vulnerable individuals through forced sex or labor. Illicit massage businesses offer a clandestine front to illicit activities by disguising themselves as legitimate businesses. This makes it challenging for law enforcement agencies and anti-trafficking organizations to detect these enterprises and their associated entities, disrupt the network, and save victims. We adopt a multi-stream data integration approach primarily focusing on consumer-generated business reviews on Yelp.com, enriched with features from contextual data sources, such as the U.S. Census and business license records. We propose a novel decision support framework that extends the traditional link prediction methods by defining a higher-order neighborhood to detect links between pairs of massage businesses and the exposure of businesses to illicit activities related to human trafficking. We achieve this by introducing a bespoke subgraph extraction strategy in GNNs where the node features are derived using NLP techniques. Comprehensive experimental results demonstrate the competitive performance of our approach over the baseline methods.

1 Introduction

The U.S. Department of State (2025) defines human trafficking as the use of force, fraud, or coercion to exploit an individual for commercial sex or labor. Human trafficking has infiltrated the massage industry due to weak regulations, governance, and laws (Organisation for Economic Cooperation and Development, 2016). Illicit massage businesses (IMBs) sell commercial sex while disguising themselves as legitimate businesses. They exploit the victims for both sex and labor while also harming the legitimate massage industry. The true extent of IMBs is unknown; however, The Network (2024) estimates that more than 15,000 IMBs operate in the U.S. These entities often operate within

a network that carries out various illicit activities, such as drug and arms trafficking, and money laundering (Miklaucic and Brewer, 2013). Given the scale of human suffering caused by their operations, several investigative agencies, both in the public and private sectors, are seeking to identify IMBs and their connections. However, identifying key indicators from the vast volume of data associated with these businesses poses a significant challenge in unearthing their discreet and intertwined operations.

Existing studies employed different approaches to identify IMBs and counter human trafficking. Aalbers and Sabat (2012); Lasker (2001); Murphy and Venkatesh (2006) applied geo-spatial analysis methods to study the spread of sexually oriented businesses, driven by the proliferation of transportation and the internet. Mletzko et al. (2018); de Vries and Radford (2021) examined the distribution of sex-trafficking offenses to determine IMB hotspots, such as highways and motels. Lugo-Graulich (2024) employed linguistic analysis to identify the distinguishing characteristics between sex-trafficking and consensual work in escort advertisements. Davy (2022) studied the role of relationships between human trafficking victims and traffickers in victim recruitment, control, and exploitation, through qualitative analysis of structured input from people with lived experiences.

With an increasing online presence, other researchers (Crotty and Bouché, 2018; Chin et al., 2019; White et al., 2021) spatially clustered IMBs and predicted their demand by analyzing customer reviews on massage boards and foot traffic data from video surveillance in combination with demographic factors. Other studies developed classification models to detect IMBs based on features extracted from customer reviews on publicly available business review websites, such as Yelp.com, either by applying sentiment analysis (Mensikova and Mattmann, 2018), or by developing a lexicon

vocabulary (Li et al., 2023). The closest study to our work (Garg et al., 2025), employed a Graph Convolutional Network (GCN) model to learn features based on the relationships between businesses, reviews, and reviewers, represented as nodes on a network. In this study, we use similar data collection and feature extraction methods. Unlike Garg et al. (2025), which classified individual massage businesses into illicit and non-illicit categories, this work focuses on the links between massage businesses.

2 Related Work and Our Contributions

2.1 Link Prediction

Link prediction aims to infer connections between entities within a domain, which may be of the same type (e.g., businesses) or different types (e.g., businesses and review writers). Early link prediction methods (Liben-Nowell and Kleinberg, 2007) leveraged similarity scores based on network measures, such as the number of common neighbors or the node degrees. The integration of embedding-based models, such as DeepWalk (Perozzi et al., 2014), Node2Vec (Grover and Leskovec, 2016), and Graph Machine Learning (GML), has provided analytical tools for processing vast amounts of data to infer links within a network. For example, GraphSAGE (Hamilton et al., 2017) aggregates neighborhood information into node embeddings and uses these embeddings for link prediction. Formulating link prediction as a classification problem, the SEAL method extracts the subgraph around the inferred link and uses a Graph Neural Network (GNN) for prediction (Zhang and Chen, 2018). Traditional GNNs assume homophily. However, the business review network studied in our work is heterophilic, where connected nodes have different labels. Specifically, businesses are connected to reviews, which, in turn, are connected to the reviewers. Neighborhood aggregation can negatively impact the performance of GNNs in such graphs. Thus, we propose a new definition of neighborhood based on higher-order links.

2.2 Link Prediction to Combat Human Trafficking

Prior studies have mainly focused on analyzing on-line escort advertisements to uncover trafficking networks. Cockbain et al. (2011) applied social network analysis to identify key hub nodes whose interdiction could significantly disrupt the network.

Szekely et al. (2015) developed a knowledge graph that includes nodes representing ads, people, locations, and contact data. Links between nodes are predicted using text and image similarity measures, as well as entity resolution methods. Chambers et al. (2019) proposed a neural network-based approach for extracting phone numbers to link escort advertisements. Li et al. (2022) combined classic rule-based and dictionary extractors with a contextualized language model to recognize entities with ambiguous names in escort ads, establishing links between these entities and the ads. Vajiac et al. (2023) proposed a micro clustering approach to identify links between escort ads. Saxena et al. (2023) predicted links between human trafficking vendors on the basis of authorship features in language patterns.

No prior work has examined business review data for link prediction between IMBs, nor has any study considered predicting different types of links. We address this gap in the literature by developing a GML-based prediction framework built on business review data. We demonstrate the applicability of the proposed framework by developing models to predict: (i) whether a pair of massage businesses are linked through their reviews and reviewers; (ii) whether the link between a pair of businesses is illicit, that is, it involves an illicit business; (iii) whether a massage business is illicit; and (iv) whether a massage business is illicit or linked (exposed) to any illicit business.

2.3 Main Contributions

Driven by investigative goals and domain expertise, the main contributions of this study include:

- Combining text embeddings extracted via Natural Language Processing (NLP) with network-based graph features and contextual features of nodes within GML to create a link prediction framework based on a business review dataset.
- Developing a subgraph extraction and labeling approach to infer higher-order links between nodes in heterophilic graphs.
- Extending traditional link prediction methods to consider different link types (i.e., illicit and non-illicit links).
- Expanding the link prediction task between two nodes to exposure detection, which aims to predict the neighborhood type of a node based on its own class and links to any other node of a certain class.

3 Methodology

3.1 Graph Construction

We represent the business review dataset as an undirected heterogeneous graph with three node types: business, review, and reviewer. Two types of edges capture the relationship between the nodes: business-review and review-reviewer edges. We refer to this graph as the *business review graph*. We focus on predicting the links between the businesses. A pair of businesses is linked if there is a path connecting them. Given the above definition of node and edge types, the minimum length of a path linking businesses A and B is 4 if the same reviewer X provided reviews for both businesses, i.e. A-reviewA-X-reviewB-B. Thus, two businesses can be linked by a path of 4, 8, 12... hops. To focus on more direct links and maximize the performance of the proposed approach, we consider predicting 4-hop and 8-hop links between business nodes in the business review graph (see Section 4.1).

3.2 Subgraph Extraction

In GML, subgraphs are used to define the receptive field or neighborhood structure that influences inference (Valsesia et al., 2023). A small receptive field may provide insufficient information, whereas larger receptive fields can lead to over-smoothing (Hamilton, 2020). In this study, we extract subgraphs from the business review graph. For prediction tasks involving pairs of businesses (link-level), we extract the m -hop ego subgraph around each business node. For tasks involving a single business (node-level), we extract the m -hop ego subgraph centered on the business node. We choose the value of the parameter m on the basis of the link definition in our experiments (see Section 4.1) as well as to balance information gain and over-smoothing.

To prevent label leakage in link prediction tasks involving business pairs, we remove all direct 4-hop links connecting the two target businesses when constructing training subgraphs. In addition, we remove 8-hop links by randomly selecting one endpoint of the link and deleting the corresponding 4-hop path connected to that endpoint. This procedure ensures that the extracted training subgraphs do not explicitly include the target links while preserving the surrounding structural context.

Each business, review, and reviewer node in the subgraph has tailored contextual features described

in Section 4.1. To capture the positional encoding of each node with respect to the target business node(s), we append z_1 and z_2 distance scores to the node features. These scores are calculated as:

$$z_1 = 1 + \min(d_A, d_B), \quad d = d_A + d_B, \quad (1)$$

$$z_2 = \begin{cases} \lfloor d/2 \rfloor (\lfloor d/2 \rfloor + d\%2 - 1) & \text{if } d < \infty, \\ 0 & \text{otherwise,} \end{cases} \quad (2)$$

where d_A and d_B are the shortest path distances from the target business nodes A and B in link prediction. For prediction tasks involving a single business node A , we only use $z_1 = 1 + d_A$.

3.3 Subgraph Encoding

We train a Graph Neural Network (GNN) to encode the structural and node features of a subgraph into an embedding vector. GNNs first propagate the information between nodes in a message-passing step. This message passing is performed over the neighborhood of each node. The next step aggregates the propagated information through a weighted sum. The final step combines the node’s features with the propagated messages from its neighborhood.

We apply message passing and aggregation over multiple neighborhood layers around each node (Kipf and Welling, 2016). The first neighborhood layer ($l = 1$) of a node consists of its immediate neighbors reachable in one hop. The second layer ($l = 2$) includes nodes reachable in two hops, and higher-order layers are defined analogously. Equation (3) demonstrates the calculation of the embedding $h_i^{(l+1)}$ for node i in layer $l + 1$.

$$h_i^{(l+1)} = \sigma \left(W_0^{(l+1)} h_i^{(l)} + \sum_{j \in N_i^{(l+1)}} W_1^{(l+1)} h_j^{(l)} \right), \quad (3)$$

where $N_i^{(l+1)}$ is the set of nodes in neighborhood layer $l + 1$ of node i , $W_0^{(l+1)}$ and $W_1^{(l+1)}$ are the self-loop and the neighborhood weight aggregation matrices, and σ is the non-linear activation function. We set the dimension of the node embedding vector h_i^l and the number of layers $l = 1, \dots, L$ using hyperparameter tuning. We maintain the same node embedding dimension in each layer, except for the last layer, where the dimension of h_i^L is set to 1. The final embedding vector (h_i) of node i is obtained by concatenating its embedding vectors in each layer, i.e. $h_i = [h_i^1, h_i^2, \dots, h_i^L]$. We sort h_i ’s based on their last entry and populate k

278 portion of them in a node embedding matrix. The
279 value of k is set using hyperparameter tuning. We
280 then process the node embedding matrix using a
281 Convolutional Neural Network (CNN) to obtain
282 an embedding vector for the entire subgraph. We
283 pass the final subgraph embedding through a fully
284 connected multi-layer perceptron for classification.

285 3.4 Benchmark Methods

286 We implement two benchmark methods for link
287 prediction between a pair of businesses: Preferen-
288 tial Attachment (PA) and Logistic Regression (LR).
289 For link class prediction or node-level prediction
290 tasks, we use LR only. PA infers a link between
291 a given pair of businesses based on the product of
292 their node degrees. LR applies a linear transforma-
293 tion to business features, followed by a nonlinear
294 activation function to produce class probabilities.

295 PA is a purely structure-based method that re-
296 quires no training but does not incorporate node
297 attributes. In contrast, LR leverages features of the
298 target business node(s) but fails to capture struc-
299 tural information from the business review graph.
300 Our proposed approach addresses these limitations
301 by jointly modeling both subgraph structural fea-
302 tures and node attributes.

303 4 Computational Experiments

304 This section describes the business review dataset,
305 the feature construction process for the three node
306 types (businesses, reviews, and reviewers), and the
307 experimental design and implementation.

308 4.1 Data

309 We leverage multi-source data to construct business
310 review graphs for massage businesses in Colorado
311 (CO), Florida (FL), and Texas (TX). We chose
312 these states for analysis because FL and TX are
313 IMB hotspots (Janis, 2020) and our collaborator,
314 [Collaborator Name], has partnerships with a lo-
315 cal law enforcement agency in CO. [Collaborator
316 Name] is a nonprofit that uses data analytics and
317 technology to fight human trafficking, which has
318 provided the following datasets:

319 **Yelp reviews:** business name, address, phone num-
320 ber, service category, price range; review text, user-
321 name, date, rating.

322 **RubMaps reviews:** business name, address, phone
323 number; review text, username, date.

324 **Business license records:** business name, address,
325 phone number, license number, license status, and

administrative orders from regulatory agencies. 326

We have augmented these datasets with contextual 327
features from publicly available data sources: 328

U.S. Census at the census tract level: demo- 329
graphic and socioeconomic variables, housing & 330
household composition, employment & industry 331
features. 332

GIS (Geographic Information System): loca- 333
tions of highways, truck stops, military bases, po- 334
lice stations, and public schools. 335

NLCD (National Land Cover Database): land 336
cover types, e.g., developed, low/high intensity. 337

Business Features. Massage businesses from the 338
Yelp dataset are geocoded to get distances to truck 339
stops, highways, military bases, police stations, 340
and schools. These places influence the location of 341
IMBs based on crime opportunity theory (de Vries, 342
2023) and stakeholder interviews (Tobey et al., 343
2022). We adopt the business labeling procedure 344
(non-illicit:0, illicit:1) from Tobey et al. (2022), 345
which utilizes the review and business features ob- 346
tained from RubMaps.ch, as well as the license 347
records. Statistically significant features are se- 348
lected using univariate logistic regression. *Table 3:* 349
Selected Data Features in Tobey et al. (Tobey et al., 350
2022) shows a complete list of business features. 351

Review Features. After applying standard natu- 352
ral language processing (NLP) techniques, such as 353
stop-word removal, lemmatization, and tokeniza- 354
tion, we convert review text into 600-dimensional 355
embedding vectors using a pretrained Doc2Vec 356
model (Li et al., 2023). These vectors are mapped 357
to a lower-dimensional space using Principal Com- 358
ponent Analysis (PCA). To create informative fea- 359
tures, we also employ the lexicon analysis from Li 360
et al. (2023), where they develop a Yelp-specific 361
lexicon for IMBs, comprising a vocabulary of 169 362
keywords selected based on their high frequency in 363
illicit reviews and input from domain experts. An 364
illicit review explicitly mentions or implies com- 365
mercial sex or other indicators of human traffick- 366
ing at a business. While commercial sex alone 367
does not constitute human trafficking unless in- 368
duced by force, fraud, or coercion U.S. Depart- 369
ment of State (2025), evidence suggests that a non- 370
negligible share of massage business workers en- 371
gaged in commercial sex are trafficking victims. 372
They weight lexicon terms by strength (1 for po- 373
tential signs for commercial sex and 2 for strong 374
indicators) and train a classifier using normalized 375
lexicon scores. We use this classifier’s output as 376

377 a review feature. We further perform sentiment
378 analysis using a RoBERTa model (Barbieri et al.,
379 2020) to classify reviews as positive, neutral, or
380 negative, and also include the review ratings (1-5)
381 as a feature.

382 **Reviewer Features.** Driven by the observa-
383 tion that IMBs predominantly serve male cus-
384 tomers (Crotty and Bouché, 2018), we use the
385 *gender guesser* package to create a gender feature
386 based on the reviewer’s username.

387 GNNs show better performance at a lower hop
388 neighborhood by avoiding over-smoothing (Hamil-
389 ton, 2020) and at higher edge homophily (Luan
390 et al., 2022) (where edge homophily defines the
391 proportion of edges that join the nodes of the same
392 class). We generalize this definition to hop ho-
393 mophily (i.e., the proportion of m -hop links that
394 connect the same class nodes). Table 2 shows that
395 4 and 8 hops exhibit high hop homophily and also
396 dense population (i.e., high total number of pairs
397 for training, testing, and validation). Therefore, we
398 select 4 and 8 hops for the link definition.

399 4.2 Design of Experiments

400 This section defines four experiments, which are
401 motivated by investigative objectives aimed at iden-
402 tifying and disrupting illicit business networks. The
403 first two experiments focus on relationships be-
404 tween pairs of businesses, whereas the last two ex-
405 periments involve individual businesses. All exper-
406 iments are formulated as binary classification tasks
407 (with Classes 1 and 0), where Class 1 consistently
408 represents businesses of investigative importance.

409 **Exp 1: Link Prediction.** This analysis enables
410 investigators to predict links between pairs of mas-
411 sages businesses. We define positive links as ob-
412 served 4-hop and 8-hop connections between busi-
413 ness node pairs (Class 1). To preserve class im-
414 balance, we sample a comparable number of busi-
415 ness node pairs that do not exhibit any connections
416 within eight hops (Class 0).

417 **Exp 2: Link Classification.** This analysis en-
418 ables investigators to prioritize the examination of
419 businesses involved in illicit connections. The ex-
420 periment predicts whether a linked pair of busi-
421 nesses includes any illicit business. We define
422 Class 1 links (illicit links) as observed 4-hop or
423 8-hop connections between business node pairs in
424 which at least one business is illicit. We sample a
425 comparable number of Class 0 links from business

426 pairs in which both businesses are non-illicit (be-
427 nign) and are connected by 4-hop or 8-hop paths.

428 **Exp 3: Business Classification.** This experiment
429 predicts whether a business is illicit (Class 1) or
430 non-illicit (Class 0). Although the primary contri-
431 bution of our work lies in link prediction and classi-
432 fication, we include this experiment to demonstrate
433 the robustness of the proposed approach in learning
434 effective embeddings for node classification.

435 **Exp 4: Illicit Exposure Detection.** This analysis
436 enables investigators to prioritize the examination
437 of businesses that are illicit or involved in illicit
438 connections. The experiment predicts whether a
439 message business is illicit or has connections to an
440 illicit business through 4 hops (Class 1). Business-
441 es in Class 0 are non-illicit and have no links to illicit
442 businesses. This experiment does not classify the
443 specific links associated with a given business. If
444 a business is predicted to belong to Class 1, the
445 analysis in Exp 2 can be applied to further identify
446 the specific illicit links.

447 4.3 Implementation

448 **Data Preparation.** We employ a three-step ap-
449 proach to generate data for model training, testing,
450 and validation in the link-based experiments (1 and
451 2). First, we perform undersampling to address
452 the class imbalance between illicit and non-illicit
453 businesses. Since non-illicit businesses are abun-
454 dant in the data, we select the ones with the most
455 reviews in our sampling approach and maintain an
456 imbalance ratio of 0.25 (Table 3). Then, we per-
457 form stratified splitting across the CO, FL, and TX
458 datasets, dividing them into 80% model develop-
459 ment and 20% held-out testing sets. We further
460 stratify the development data into 80% training and
461 20% validation sets, with the validation set used
462 for early stopping. This inductive splitting across
463 businesses prevents data leakage and aligns with
464 the purpose of our framework, which is to make in-
465 ferences about unknown businesses and networks.
466 Finally, we sample training, validation, and testing
467 business pairs from their respective business-review
468 graphs (Table 4) using the class definitions in each
469 experiment. In particular, we sample up to 10,000
470 business pairs for Class 1 and up to 10,000 pairs for
471 Class 0, considering all pairs if the set is smaller.
472 The data preparation for Experiments 3 and 4 in-
473 volves only the stratified splitting of businesses into
474 training, validation, and testing sets.

Hyperparameter Tuning. The key hyperparameters chosen for tuning along with their search space include: number of GNN layers $L \in \{2, 4, 6\}$, hidden dimension of the node embeddings for the GNN layers $|h_i^L| \in \{8, 16, 32\}$, node pool percentile $k \in \{0.3, 0.6, 0.9\}$ as well as the number of input and output channels of the convolutional layer $(c_1, c_2) \in \{(4, 8), (8, 16)\}$.

We use the FL dataset for hyperparameter optimization, as it is the medium-sized set in terms of number of businesses, reviews, and reviewers (Table 5). We perform separate hyperparameter tuning for each experiment based on the AUC value and use the finalized hyperparameters (as shown in Table 6) for a single-run testing. For Experiments 1 and 2, we sample 2,500 business pairs in Class 0 and Class 1 from the businesses in the training set. For Experiments 3 and 4, we consider the entire training set, as these experiments are based on businesses rather than business pairs. We further perform 5-fold stratified cross-validation across all experiments and report mean results across all folds with their sample standard deviations.

We report two performance metrics, AUC and Average Precision (Avg Prec), i.e., the area under the precision-recall curve. These threshold-agnostic metrics are particularly suitable for imbalanced datasets. We implement the framework on Google Colab (virtualized Linux) using Python 3.11, equipped with an NVIDIA L4 GPU and an Intel Xeon 2.20 GHz CPU.

4.4 Numerical Results

Table 1 presents the performance of the GML-based prediction model (GNN) as well as the performance of the LR and PA methods over the held-out testing sets for CO, FL, and TX.

Link-level Experiments. The GML-based methods consistently achieve the highest AUC and Avg Prec across CO and TX in Experiments 1 (Link Prediction) and 2 (Link Classification), showcasing robust performance even with the added difficulty of classification on top of structural link prediction. We see a notable improvement of 0.2692 and 0.2909 in the AUC and Avg Prec for CO (dataset with the lowest labeled businesses) compared to the LR model, highlighting the importance of high-order links and network-informed learning in low-data regime. The LR model, which is trained on handcrafted business features, also outperforms PA, with the effect more pronounced in larger datasets

(FL and TX), demonstrating the importance of features and model training over untrained network measures. Competitive LR results for FL across both experiments suggest that the subgraph structure is homogeneous for both classes.

Node-level Experiments. In Experiment 3 (Node Classification), we observe significant improvements in LR performance metrics, suggesting that direct business features are highly predictive. Even though LR dominates, the learned embeddings from the GNN model show competitive performance and do not degrade. GNN achieves an extremely high Avg Prec score of 0.9920 for TX in Experiment 4 (Illicit Exposure Detection) with substantial improvements over LR for CO and FL, establishing the criticality of neighborhood aggregation through network-informed learning in detecting exposure.

Class Imbalance. We test the performance of the GNN, LR and PA models at various imbalance ratios to ascertain generalizability for link-level experiments. In Figure 1, for CO, while LR is the least sensitive, GNN outperforms the former across all ratios for both Exp 1 and 2.

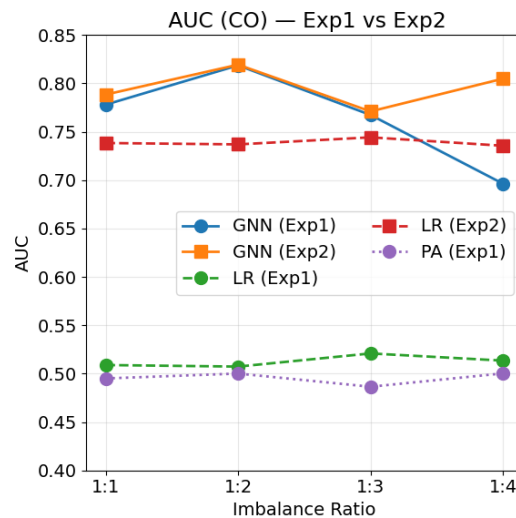


Figure 1: Model performance at different imbalance ratios for link-level experiments. Analogous results for FL and TX are presented in Figure 6 and Figure 7.

4.5 Parameter Sensitivity

Number of GNN Layers (L). We vary the number of GNN layers within $\{2, 4, 6\}$ while fixing other parameters to assess the implication of neighborhood layers towards nodes' embeddings. Figure 2 suggests that the model's AUC is highly sen-

Exp	Method	CO		FL		TX	
		AUC	Avg Prec	AUC	Avg Prec	AUC	Avg Prec
1	GNN	0.7781	0.8076	0.6492	0.6984	0.7324	0.7876
	LR	0.5089	0.5167	0.6468	0.6068	0.6015	0.5722
	PA*	0.4951	0.4888	0.5202	0.5225	0.4873	0.5020
2	GNN	0.7885	0.7718	0.7725	0.8056	0.5485	0.5304
	LR	0.7384	0.7666	0.7733	0.8200	0.5000	0.5000
3	GNN	0.6626	0.5912	0.8971	0.8265	0.8580	0.7726
	LR	0.8097	0.5093	0.9242	0.8968	0.9596	0.9138
4	GNN	0.5950	0.9387	0.7414	0.9434	0.8534	0.9920
	LR	0.5029	0.9014	0.6145	0.9136	0.6530	0.9780

Table 1: Prediction performance across four experiments. * PA is only used in Experiment 1 because it is suitable for link prediction, not for link classification.

sitive in Exp 2, showing an improvement of 0.2263 when the number of layers is increased from 2 to 6. This reinforces the impact of GNN’s message passing and aggregation towards creating informative representations for predicting the link class. In contrast, for the simplistic task of link prediction in Exp 1, we see a decrease in performance, which is consistent with the phenomenon of over-smoothing. Exp 3 (Node Classification) shows the lowest sensitivity, while performance improves as L increases for Exp 2 and 4 (Link Classification and Exposure Detection), due to their reliance on neighborhood context.

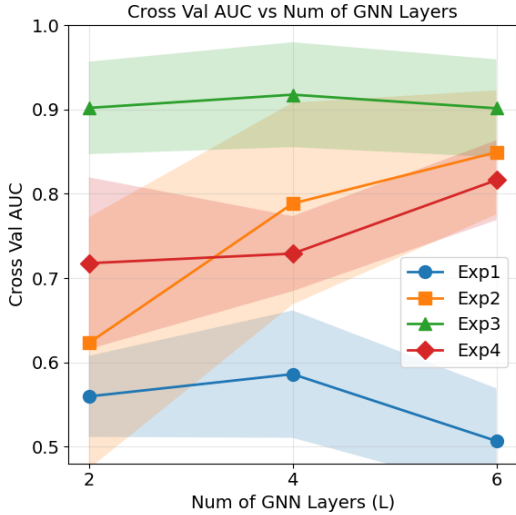


Figure 2: Model cross-validation AUC \pm std dev values with different numbers of layers on the FL dataset.

Node Pool Percentile (k). We examine the node pool size by varying k within $\{0.3, 0.6, 0.9\}$. A k value of 0.6 is interpreted as the 60th percentile among all training subgraph sizes (number of

nodes) to be chosen as the pooling size. In Figure 3, Exp 1, 2, and 4 show peak AUC at $k = 0.6$, which prevents over-smoothing due to noise from weakly informative nodes. Exp 3 shows minimal sensitivity to k , indicating that business classification depends primarily on its own features.

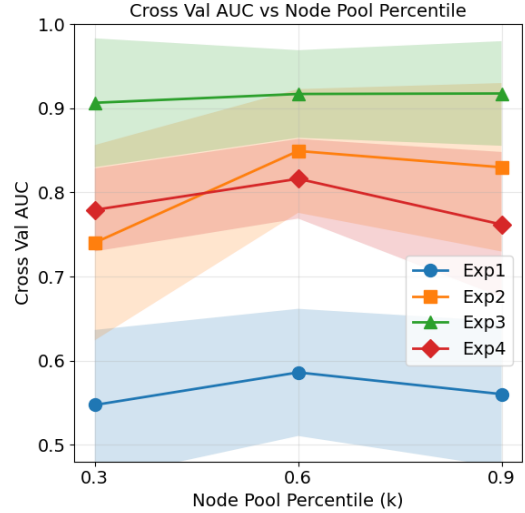


Figure 3: Model cross-validation AUC \pm std dev values with different node pool percentile on the FL dataset.

4.6 Ablation Study

In this section, we analyze the contribution of the z_1 and z_2 positional encoding scores in the performance of the GNN model. To do so, we define a new model named GNN_woz, i.e., GNN without these scores. Table 7 shows that GNN outperforms GNN_woz at most instances of link-level experiments, establishing the value of network topology. For the node-based experiments, where node features dominate (as evidenced by the improved performance of LR), both models exhibit comparable

performance.

5 Feature Importance

To assess the importance of each feature in the GNN model, we leverage a gradient-based method called Integrated Gradients (IG) (Sundararajan et al., 2017), which we implement using Captum (Kokhlikyan et al., 2020). IG evaluates the model at m different inputs created by linearly interpolating the original input features x and a neutral baseline x' . It then averages the gradients of the output with respect to those inputs to generate an attribution score IG_i , which captures the importance of the i_{th} feature, defined as:

$$IG_i = (x_i - x'_i) \cdot \frac{1}{m} \sum_{k=1}^m \frac{\partial F(x' + \frac{k}{m}(x - x'))}{\partial x_i}, \quad (4)$$

We calculate the average of the absolute values of these scores across all nodes and subgraphs in the testing set to generate final feature-level attributions. For this, we use the GNN_woz model to focus on business features. Figure 4 depicts the top five most important business features in the test set for Exp 2 (Link Classification) for CO. To generate attribution scores for a link, we can calculate the sum of the absolute values of the attribution scores of the endpoint nodes and consider this as a proxy for link importance. Figure 8 depicts a subgraph with its links highlighted on the basis of these attributions.

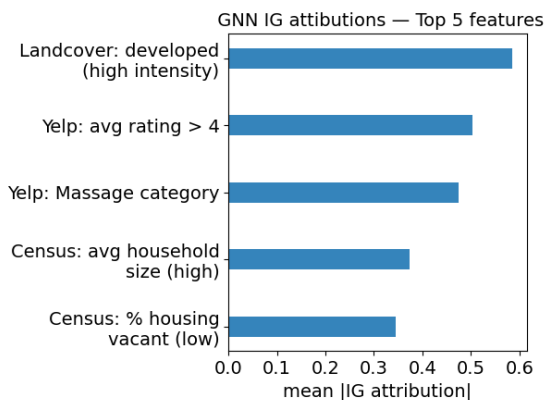


Figure 4: Mean IG attributions from the GNN Model.

To compare the feature importance of GNN_woz with LR, we generate the SHAP values (Lundberg and Lee, 2017) of the LR features. Figure 5 presents a beeswarm plot highlighting the five most influential features based on SHAP values.

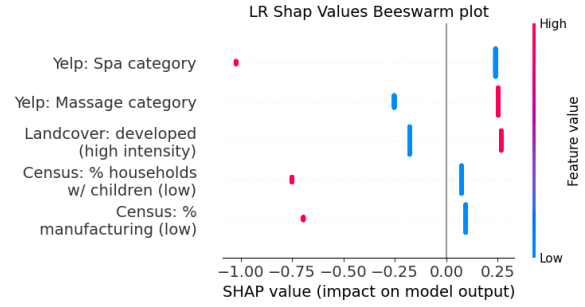


Figure 5: Beeswarm plot of the shap values.

In Figures 4 and 5, two of the top five features (explained in Table 8) appear in both models, showcasing consistency in identifying salient signals. Since the GNN model leverages the network structure and message passing across business, review, and reviewer nodes, it identifies different business features of high importance for classifying a link.

6 Discussion and Conclusions

This work studies the network between illicit massage businesses using multifaceted business-review data and extends the link prediction models to classify links and detect exposure. The LR model utilizes handcrafted business features, leveraging Natural Language Processing and the domain expertise of collaborators, and builds upon a logistic regression classifier without incorporating network information. To overcome this, we introduce a graph machine learning-based GNN model that learns latent structural relationships between business, review, and reviewer nodes in a business review dataset and infers high-order links between businesses. We introduce four experiments motivated by investigative goals. Two link-level experiments, which aim to predict and classify the link between two business nodes, and two node-level experiments, which classify a business node and detect its exposure to illicit business nodes. The GNN model outperforms LR in the former set of experiments, where the relative structure of the subgraph around the target business nodes is key, whereas LR dominates when individual business features are crucial. In this work, we focus on learning a network to infer links between business nodes; an extension of this work can incorporate predicting links between different node types. Additionally, as we observed improved performance in both the smallest and largest datasets, future work will aim to study characteristic measures of the graph topology to elucidate this performance gain.

7 Ethical Considerations

The proposed framework can enhance the transparency of disciplinary actions and inform regulatory policies that protect vulnerable industries, such as the massage industry. Furthermore, with comprehensive experiments driven by investigative goals to identify illicit massage businesses and their connections, the models demonstrate a real-world use case that can serve as a decision-support tool for law enforcement agencies and counter-human trafficking organizations, facilitating the proper allocation of investigative resources to uncover a hidden network of illicit businesses. The classification models proposed in this work can generate false positive results. Therefore, the models should be used cautiously to inform decision-making and prioritize investigations. This work supports reproducibility without raising ethical concerns or posing risks to society. Additionally, the datasets were constructed in accordance with strict ethical guidelines, ensuring the anonymity of businesses and reviewers. The reviews used in this work are redacted to remove the names of individual people. However, it can contain offensive content about commercial sex. To protect data, we share it only with relevant researchers and investigators through a data use agreement.

8 Limitations

We employ the inductive splitting of businesses into training, validation, and testing sets to prevent data leakage; however, this results in the training graphs being larger in size than the validation and testing graphs. A nuanced splitting strategy to create disjoint business sets with graphs of similar characteristics can improve model performance. We also limit sampling to 2,500 business pairs for hyperparameter optimization and 10,000 business pairs for testing to control computational complexity. Multiple replications of hyperparameter tuning and model testing with additional samples will enhance the robustness of our results. We used three datasets containing Yelp reviews for massage businesses in CO, FL, and TX, which limits the generalizability of the results to other states with a significant domain shift; however, the proposed model is generalizable to other datasets. Finally, manual labeling of the data required domain knowledge and careful reasoning, which constrained the availability of annotated data.

709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764

References

Manuel B. Aalbers and Magdalena Sabat. 2012. [Re-making a Landscape of Prostitution: the Amsterdam Red Light District](#). *City*, 16(1-2):112–128.

Francesco Barbieri, José Camacho-Collados, Leonardo Neves, and Luis E. Anke. 2020. [Tweeteval: Unified benchmark and comparative evaluation for tweet classification](#). *CoRR*, abs/2010.12421.

Nathanael Chambers, Timothy Forman, Catherine Griswold, Kevin Lu, Yogaish Khastgir, and Stephen Steckler. 2019. [Character-based models for adversarial phone extraction: Preventing human sex trafficking](#). In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 48–56.

John J. Chin, Lois M. Takahashi, and Douglas J. Wiebe. 2019. [Where and Why Do Illicit Businesses Cluster? Comparing Sexually Oriented Massage Parlors in Los Angeles County and New York City](#). *Journal of Planning Education and Research*, 43(1):106–121.

Eleanor Cockbain, Helen Brayley, and Gloria Laycock. 2011. [Exploring internal child sex trafficking networks using social network analysis](#). *Policing: A Journal of Policy and Practice*, 5(2):144–157.

Sean M. Crotty and Vanessa Bouché. 2018. [The Red-Light Network: Exploring the Locational Strategies of Illicit Massage Businesses in Houston, Texas](#). *Papers in Applied Geography*, 4(2):205–227.

Deanna Davy. 2022. [Trafficked by someone i know: A qualitative study of the relationships between trafficking victims and human traffickers in albania](#). Research report, UNICEF Albania & IDRA. Qualitative study using semi-structured interviews with trafficking survivors and key informants.

Ieke de Vries. 2023. [Examining the geography of illicit massage businesses hosting commercial sex and sex trafficking in the united states: The role of census tract and city-level factors](#). *Crime & Delinquency*, 69(11):2218–2242.

Ieke de Vries and Jason Radford. 2021. [Identifying online risk markers of hard-to-observe crimes through semi-inductive triangulation: The case of human trafficking in the United States](#). *The British Journal of Criminology*, 62(3):639–658.

Vasuki Garg, Osman Y. Özaltın, Maria E. Mayorga, and Sherrie Bosisto. 2025. [Detecting illicit massage businesses by leveraging graph machine learning](#). In *Proceedings of the Thirty-Fourth International Joint Conference on Artificial Intelligence, IJCAI-25*, pages 9647–9655. International Joint Conferences on Artificial Intelligence Organization. AI and Social Good.

Aditya Grover and Jure Leskovec. 2016. [node2vec: Scalable feature learning for networks](#). In *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 855–864. ACM.

William L. Hamilton. 2020. [Graph representation learning](#). *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 14(3):1–159. 765
766
767

William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. [Inductive representation learning on large graphs](#). In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 30, pages 1024–1034. 768
769
770
771

Elizabeth Ranade Janis. 2020. [Unmasking Trafficking in Illicit Massage Businesses Across the United States - Human Trafficking Institute](#) — traffickinstitute.org. <https://traffickinginstitute.org/illicit-message-businesses/>. [Accessed 11-06-2024]. 772
773
774
775
776
777

Thomas N. Kipf and Max Welling. 2016. [Semi-supervised classification with graph convolutional networks](#). *CoRR*, abs/1609.02907. 778
779
780

Narine Kokhlikyan, Vivek Miglani, Miguel Martin, Edward Wang, Bilal Alsallakh, Jonathan Reynolds, Alexander Melnikov, Natalia Kliushkina, Carlos Araya, Siqi Yan, and Orion Reblitz-Richardson. 2020. [Captum: A unified and generic model interpretability library for pytorch](#). *CoRR*, abs/2009.07896. 781
782
783
784
785
786

Stephanie Lasker. 2001. [Sex and the city: zoning pornography peddlers and live nude shows](#). *UCLA Law Review*, 49:1139–1185. 787
788
789

Ruoting Li, Margaret Tobey, Maria E. Mayorga, Sherrie Caltagirone, and Osman Y. Özaltın. 2023. [Detecting Human Trafficking: Automated Classification of Online Customer Reviews of Massage Businesses](#). *Manufacturing & Service Operations Management*, 25(3):1051–1065. 790
791
792
793
794
795

Yifei Li, Pratheeksha Nair, Kellin Pelrine, and Reihaneh Rabbany. 2022. [Extracting person names from user generated text: Named-entity recognition for combating human trafficking](#). In *Findings of the Association for Computational Linguistics: ACL 2022*, pages 2854–2868, Dublin, Ireland. Association for Computational Linguistics. 796
797
798
799
800
801
802

David Liben-Nowell and Jon Kleinberg. 2007. [The link-prediction problem for social networks](#). *Journal of the American Society for Information Science and Technology*, 58(7):1019–1031. 803
804
805
806

Sitao Luan, Chenqing Hua, Qincheng Lu, Jiaqi Zhu, Mingde Zhao, Shuyuan Zhang, Xiao-Wen Chang, and Doina Precup. 2022. [Revisiting heterophily for graph neural networks](#). *Preprint*, arXiv:2210.07606. 807
808
809
810

Kristina Lugo-Graulich. 2024. [Indicators of Sex Trafficking in Online Escort Ads, 7 U.S. States, 2013–2020](#). Inter-university Consortium for Political and Social Research, Ann Arbor, MI. ICPSR distributor; dataset. 811
812
813
814
815

Scott M Lundberg and Su-In Lee. 2017. [A unified approach to interpreting model predictions](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc. 816
817
818
819

820 Anastasija Mensikova and Chris A. Mattmann. 2018. Ensemble Sentiment Analysis to Identify Human
821 Trafficking in Web Data. In *Workshop on Graph
822 Techniques for Adversarial Activity Analytics (GTA
823 2018)*, Marina Del Rey, CA, USA, pages 5–9. 873

825 Michael Miklaucic and Jacqueline Brewer, editors. 2013. *Convergence: Illicit Networks and National Security
826 in the Age of Globalization*. National Defense Uni- 874
827 versity Press, Washington, DC. 875

829 Deborah Mletzko, Lucia Summers, and Ashley N. Arnio. 2018. Spatial patterns of urban sex trafficking. *Journal of Criminal Justice*, 58:87–96. 876

832 Alexandra K. Murphy and Sudhir A. Venkatesh. 2006. Vice Careers: The Changing Contours of Sex Work
833 in New York City. *Qualitative Sociology*, 29(2):129–
834 154. 877

836 Organisation for Economic Co-operation and Develop- 878
837 ment. 2016. *Trafficking in Persons and Corruption:
838 Breaking the Chain*. OECD Publishing, Paris. Chap- 879
839 ter 1: Trafficking in persons: Weak governance and 880
840 growing profits. 881

841 Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. 2014. Deepwalk: Online learning of social representations.
842 In *KDD*. 882

844 Vageesh Saxena, Benjamin Bashpole, Gijs Van Dijck, and Gerasimos Spanakis. 2023. Idtraffickers: An
845 authorship attribution dataset to link and connect po- 883
846 tential human-trafficking operations on text escort 884
847 advertisements. *arXiv preprint arXiv:2310.05484*. 885

849 Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. *CoRR*,
850 abs/1703.01365. 886

852 Pedro Szekely, Craig A. Knoblock, Jason Slepicka, Andrew Philpot, Amandeep Singh, Chengye Yin,
853 Dipsy Kapoor, Prem Natarajan, Daniel Marcu, Kevin 887
854 Knight, David Stallard, Subessware S. Karunamoorthy, Rajagopal Bojanapalli, Steven Minton, Brian
855 Amanatullah, Todd Hughes, Mike Tamayo, David 888
856 Flynt, Rachel Artiss, and 4 others. 2015. Building 889
857 and using a knowledge graph to combat human traf- 890
858 ficking. In *Proceedings of the 14th International
859 Semantic Web Conference (ISWC)*, Lecture Notes in
860 Computer Science, Cham. Springer. 891

863 The Network. 2024. What is the Illicit Massage Industry? | The Network — thenetworkteam.org.
864 [https://www.thenetworkteam.org/research/
865 what-is-the-illicit-massage-industry](https://www.thenetworkteam.org/research/what-is-the-illicit-massage-industry). 892
866 [Accessed 08-02-2025]. 893

868 Margaret Tobey, Ruoting Li, Osman Y. Özaltın, Maria E. Mayorga, and Sherrie Caltagirone. 2022. Interpret-
869 able models for the automated detection of hu- 894
870 man trafficking in illicit massage businesses. *IISE
871 Transactions*, 56(3):311–324. 895

U.S. Department of State. 2025. Understanding Human Trafficking. [https://www.state.gov/
872 what-is-trafficking-in-persons/](https://www.state.gov/what-is-trafficking-in-persons/). [Accessed
873 05-02-2025]. 874

Catalina Vajiac, Meng-Chieh Lee, Aayushi Kulshrestha, Sacha Levy, Namyong Park, Andreas Olligschlaeger,
875 Cara Jones, Reihaneh Rabbany, and Christos Faloutsos. 2023. Deltashield: Information theory for
876 human- trafficking detection. *ACM Trans. Knowl.
877 Discov. Data*, 17(2). 878

Diego Valsesia, Giulia Fracastoro, and Enrico Magli. 2023. Ran-gnns: Breaking the Capacity Limits
879 of Graph Neural Networks. *IEEE Transactions on
880 Neural Networks and Learning Systems*, 34(8):4610–
881 4619. 882

Anna White, Seth Guikema, and Bridgette Carr. 2021. Why are You Here? Modeling Illicit Massage Busi-
883 ness Location Characteristics with Machine Learning. 884
885 *Journal of Human Trafficking*, 10(1):20–40. 886

Muhan Zhang and Yixin Chen. 2018. Link prediction based on graph neural networks. *CoRR*,
887 abs/1802.09691. 888

A Appendix: Business Pairs and Homophily

Data	Hop	Connected Business Pairs	Hop Homophily
CO	4	3715	81.5%
	8	40326	79.2%
	12	37335	65.4%
	16	2837	23.6%
FL	20	42	90.5%
	4	9171	87.1%
	8	134336	85.2%
	12	122394	64.2%
TX	16	12386	21.3%
	20	1052	39.6%
	4	22121	84.4%
	8	413551	81.1%
TX	12	268967	55.0%
	16	14416	32.2%
	20	342	71.6%

Table 2: Number of connected business pairs and homophily across hops for CO, FL, and TX.

B Appendix: Undersampled Business Data Statistics

Dataset	Total	Illicit	Non-Illicit	Ratio
CO	425	85	340	0.25
FL	785	157	628	0.25
TX	1230	246	984	0.25

Table 3: Business statistics for CO, FL, and TX.

C Appendix: Class Distribution across the Datasets after Inductive Splitting

Exp	Dataset	Split	Class 1	Class 0
1	CO	Train	13311	23545
		Val	305	1973
		Test	765	2805
		Total	14381	28323
	FL	Train	44790	80961
		Val	975	6900
		Test	1792	10454
		Total	47557	98315
	TX	Train	141470	167821
		Val	3291	16015
		Test	6747	23388
		Total	151508	207224
2	CO	Train	2715	10596
		Val	83	222
		Test	214	551
		Total	3012	11369
	FL	Train	6614	38176
		Val	111	864
		Test	204	1588
		Total	6929	40628
	TX	Train	25483	115987
		Val	543	2748
		Test	1564	5183
		Total	27590	123918

Table 4: Total Class 1 and Class 0 links across splits for Experiments 1 and 2.

D Appendix: Node Types and Counts

Data	Business	Review	Reviewer
CO	425	7662	5523
FL	785	13584	9335
TX	1230	21824	13508

Table 5: Number of businesses, reviews, and reviewers in CO, FL, and TX.

E Appendix: Hyper-parameters used in the Four Experiments

Parameter	Exp 1	Exp 2	Exp 3	Exp 4
L	4	6	4	6
$ h_i^l $	8	32	32	32
k	0.6	0.6	0.9	0.6
(c_1, c_2)	(4,8)	(4,8)	(4,8)	(4,8)
m	6	6	6	6
Loss Function	BCE	BCE	BCE	BCE
Optimizer	Adam	Adam	Adam	Adam
Max Epochs	100	100	100	100

Table 6: Key hyperparameters across experiments. The first four parameters are tuned. BCE: Binary Cross Entropy.

F Appendix: Ablation Study

Dataset	Exp	Method	AUC	Avg Prec
CO	1	GNN	0.7781	0.8076
		GNN_woz	0.7030	0.6921
	2	GNN	0.7885	0.7718
		GNN_woz	0.5022	0.5049
	3	GNN	0.6626	0.5912
		GNN_woz	0.6280	0.5260
	4	GNN	0.5950	0.9387
		GNN_woz	0.5950	0.9404
FL	1	GNN	0.6492	0.6984
		GNN_woz	0.6059	0.5886
	2	GNN	0.7725	0.8056
		GNN_woz	0.4963	0.4891
	3	GNN	0.8971	0.8265
		GNN_woz	0.9188	0.8682
	4	GNN	0.7414	0.9434
		GNN_woz	0.7091	0.9343
TX	1	GNN	0.7324	0.7876
		GNN_woz	0.6233	0.6111
	2	GNN	0.5485	0.5304
		GNN_woz	0.5733	0.5748
	3	GNN	0.8580	0.7726
		GNN_woz	0.8510	0.7129
	4	GNN	0.8534	0.9920
		GNN_woz	0.8197	0.9892

Table 7: Prediction performance of the GNN model with and without (z_1, z_2) positional encodings across datasets.

905
906
907

G Appendix: Model Performance at different Imbalance ratios for Link-level Experiments for FL and TX

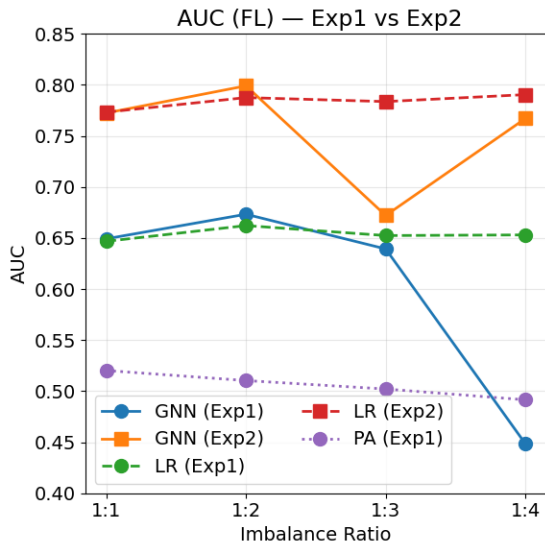


Figure 6: Model performance at different imbalance ratios for link-level experiments for FL.

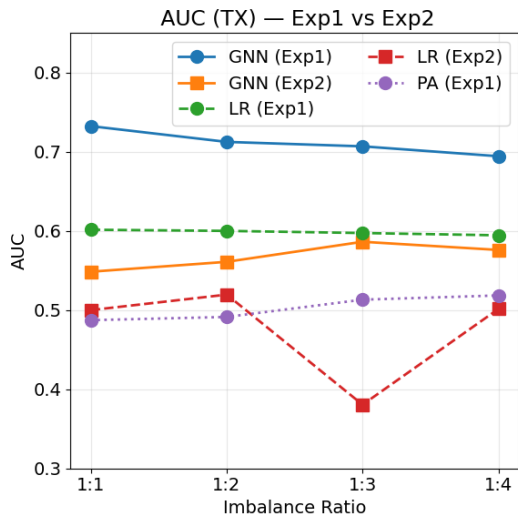


Figure 7: Model performance at different imbalance ratios for link-level experiments for TX.

H Appendix: Subgraph with Link Importance

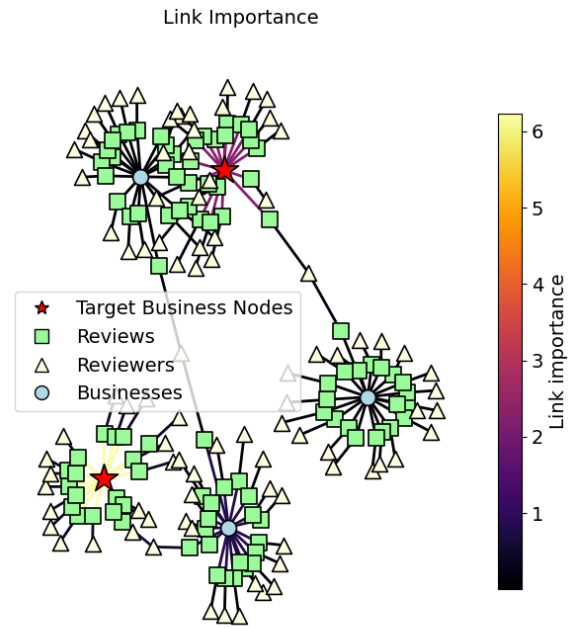


Figure 8: A Class 1 (Exp 2) subgraph with highlighted link importance.

I Appendix: Description of Business Features

For the complete list, please refer to *Table 3: Selected Data Features* in (Tobey et al., 2022)

910
911
912
913

908
909

Features	Description
Yelp: Spa category	if business categorizes under Day Spas, Medical Spas, Saunas, Float-Spa, Beauty & Spas
Yelp: Massage category	if business categorizes under Massage Therapy, Massage, Reflexology, Reiki, Tui-Na & Massage Schools
Yelp: average rating > 4	if the average of all Yelp review rating is greater than 4
Census: % households with children (low)	if percentage of households with children in the zip code is low
Census: % manufacturing industry (low)	if percentage of people employed in the zip code are in the manufacturing industry is low
Census: % housing vacant (low)	if percentage of vacant housing in the zip code is low
Census: average household size (high)	if the average household size in the zip code is low
Landcover: developed (high intensity)	if landcover type in the zip is under developed high intensity (NLCD)

Table 8: Description of binary business features.