

# Fundamentals of Data Science

Martin Aviles, Jeremy Palmerio, Mariëtte Olijslager, Andrei Maria

University of Amsterdam

**Abstract.** Link prediction is a widely-used technique for forecasting potential connections between nodes in a network, leveraging existing relationships and engineered features to make accurate predictions. Throughout the scope of the paper, a network is used to train a model to predict future links in that network. Each node possesses one of six 'sensitive' features, which provide some insight into whether connections exist between some other nodes. Given those features, various other features were engineered using properties of the network, such as the Jaccard Index, degree, triadic score, community score, and the Soundarajan scores. Using the engineered and base features, 3 prediction techniques were employed using machine learning techniques, namely a Logistic Regression model, Random Forrest, and an XGBoost model. Through parameter tuning, the XGBoost model was the highest performing model, achieving an AUC score of 0.94. The fairness of the model was then enhance, using threshold optimization to increase fairness metrics such as demographic parity, and equalized odds. Despite the resulting fairness improvement, the overall performance of the model dropped significantly. The second part of the paper evaluates the social and ethical concerns regarding the use of fairness in solving problems within similar spaces. The ethical section accounts for the roles and responsibilities of large entities in reconsidering their practices in managing their data, with respect to fairness, privacy and transparency. A primary focus was placed on fairness, with respect to the TikTok algorithms, and lawsuits around the algorithm, positioned toward the targeting of children and vulnerable groups to maximise retention on the application.

## 1 Technical solution

### 1.1 Introduction

Link prediction is an area of data science in which connections within large networks can be inferred through feature engineering, and the application of various machine learning techniques. It is a common method for recommendation algorithms [5]. The aim of this paper is to take a given network, with an arbitrary set of 'sensitive' features, and predict whether links exist between nodes in the network. The dataset given for the purpose of research in this paper included 1500 nodes, each possessing one of six features, and 6600 network edges. The implication of conducting this study, for instance, extends to solving problems within the areas of security and social networks, knowledge graphs, transportation, and fraud detection. On one hand the objective of this paper is to develop, train and optimize a classification model to predict future links in the network. On the other hand, we aim to explore the ethical aspect of designing and deploying such a model, by delving into a TikTok case study.

### 1.2 Data Preprocessing

For this project, the raw data was provided as a CSV file consisting of a list of nodes and their respective attributes. Each node  $n_i$  is represented as a tuple, where the first element is its ID, and its second element is its respective attribute.

In addition, a list of edges is provided in an `.edgelist` file, where each tuple represents a connection between two nodes. An edge between nodes  $n_i$  and  $n_j$  is denoted as:

$$e_{ij} = (n_i, n_j)$$

where:

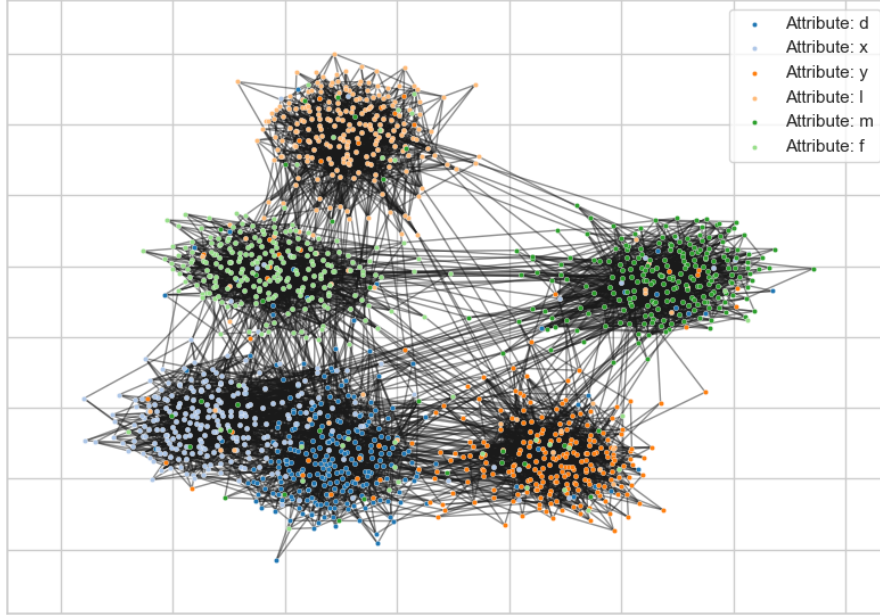
$$e_{ij} \in E$$

represents the set of all edges  $E$ , and each  $e_{ij}$  indicates a connection between node  $n_i$  and node  $n_j$ .

The python package `networkx` was employed in creating a graph object of the given nodes and edges, due to its capabilities in visualizing, manipulating, and retrieving specific pieces of information about nodes and their connections. Furthermore, this package is used in conjunction with `pandas` and `numpy` in order to perform table operations easily and efficiently, while allowing compatibility with machine learning models.

### 1.3 Data Visualization

A visualization was made to identify any underlying patterns within the data. We used the spring layout algorithm from `networkx` in order to compute the positions of the nodes. This layout simulates a force-directed graph, where nodes are positioned based on a physical model, pulling related nodes closer together, while non-connected nodes repel each other. Using this algorithm, communities between groups of nodes are easily identified.



**Fig. 1.** Representation of the network with nodes colored according to their attributes. The network layout is generated with the `networkx spring_layout`.

From the above graph, it was evident that nodes with certain attributes were closely related in communities, therefore the node attributes appeared to be an important factor in determining links between given nodes. From the above graph, there also appear to be outliers within those communities, distinguished by minority colors standing out from the majority. This suggests that attributes alone may not be sufficient alone in predicting links, but rather other features and deciding attributes must then be derived from the current dataset to account for the outliers within the communities shown.

### 1.4 Feature engineering

In order to predict the existence or not of a link between a given pair of nodes  $(n_i, n_j)$ , features are calculated. Two different classes of features are considered in this paper: similarity, structural. The final subset of features was refined iteratively by selecting the features with the highest importance in the Random Forest model.

**Structural Features** These refer to characteristics that describe the arrangement or relationships between nodes within a network or graph. They focus on the topology or structure of the network. The following features were calculated:

- Degree: These represent the degrees (number of edges) of nodes  $i$  and  $j$  in the graph, respectively. The degree of a node indicates how many connections it has.

$$\text{degree}(i) = \sum_{j \in N(i)} 1$$

- Triadic score: This measures the number of triangles that nodes  $i$  and  $j$  form with their common neighbors. It's a form of clustering that checks if two connected nodes also have a common third node.

$$T(i, j) = \frac{|\{k : (i, k) \in E, (j, k) \in E\}|}{\text{max possible triangles}}$$

- Community Score: This measures whether nodes  $i$  and  $j$  belong to the same community.

$$COM(i, j) = \begin{cases} 1 & \text{if } C(i) = C(j) \\ 0 & \text{if } C(i) \neq C(j) \end{cases}$$

- Soundarajan Score: This feature checks the number of common neighbors that  $i$  and  $j$  share and belong to the same community.

$$S(i, j) = \sum_{k \in N(i) \cap N(j)} COM(i, k) \cdot COM(j, k)$$

**Similarity Features** These focus on how similar two nodes are, often based on shared attributes or connections. In this task, the following metrics were included in the feature space:

- Common Neighbors: Counts the number of shared neighbors between nodes  $i$  and  $j$ .

$$CN(i, j) = |N(i) \cap N(j)|$$

- Jaccard Score: Measures the similarity between nodes  $i$  and  $j$  based on their shared neighbors relative to their total neighbors.

$$J(i, j) = \frac{|N(i) \cap N(j)|}{|N(i) \cup N(j)|}$$

- Adamic-Adar Score: Gives more weight to common neighbors that have fewer total connections.

$$AA(i, j) = \sum_{k \in N(i) \cap N(j)} \frac{1}{\log(\text{degree}(k))}$$

- Preferential Attachment: Measures the likelihood of nodes  $i$  and  $j$  connecting based on the product of their degrees.

$$PA(i, j) = \text{degree}(i) \cdot \text{degree}(j)$$

- Cosine Similarity: Measures the similarity between nodes  $i$  and  $j$  based on their attributes.

$$\text{Cosine}(i, j) = \frac{\mathbf{A}_i \cdot \mathbf{A}_j}{\|\mathbf{A}_i\| \cdot \|\mathbf{A}_j\|}$$

where  $\mathbf{A}_i$  and  $\mathbf{A}_j$  are attribute vectors of nodes  $i$  and  $j$ .

## 1.5 Data preparation

After the preprocessing and general data cleaning were performed and predictive and response sets were created. First, a training graph was generated by removing 20% of the edges from the original graph. These removed edges were later used as a testing set to evaluate the models. This graph split was performed to avoid leaking data from features that use the whole graph. Given a list of node pairs, each feature was calculated and makes up the total feature space; these make up the positive examples (nodes we know connect). In addition, negative examples were generated by randomly sampling the set of node pairs that are not connected. Iteratively, a 1:1 ratio of positive to negative samples was chosen to avoid creating imbalance in the data. The response column was also generated with 1 and 0 for positive and negative examples, respectively.

## 1.6 Model selection

For link prediction, three models were investigated: a Logistic Regression, a Random Forest and an XGBoost model. Implemented with `sklearn`, the Logistic Regression was aimed at providing a baseline reading, ignoring potential non-linearities and allowing for better interpretation. It has been showed to perform well on binary link prediction [4]. Also using `sklearn`, the Random Forest was used to capture more complex relationships including non-linearities. It is also more robust to noise and allows for feature importance to be calculated, which can help against over-fitting [4]. Finally, the XGBoost model was implemented using the XGBoost Python package and was chosen for its robust predictive and computational performance. Indeed, its ability to handle complex feature interaction has been put forward in the literature [7], while also beating most competition in terms of computing time. However, while the logistic model lacks the incorporation of non-linearities, the two tree-based models require careful grid searching to ensure proper regularization to avoid over-fitting.

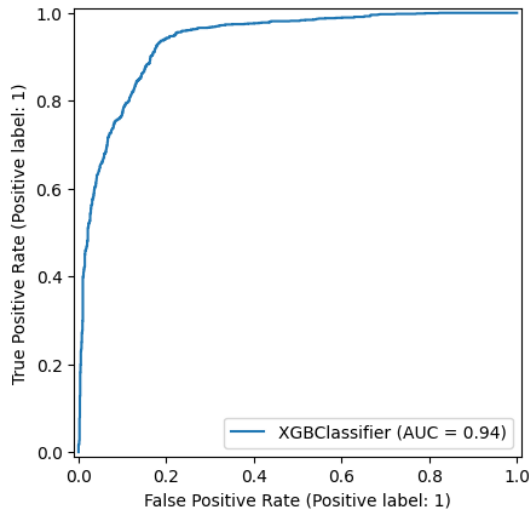
## 1.7 Model training and validation

For each model, a parameter grid was defined based on common value ranges (see Tables 4, 5, 6 in Appendix A). Then, a grid search was performed using `sklearn`'s `GridSearchCV`, which allows for cross-validation to tune the hyperparameters. A 5-fold cross-validation was chosen based on common practice. The best models were then evaluated by comparing their train and test accuracies, confusion matrices, and AUC scores. Finally, the best model was trained on the whole full set of 1500 nodes and 6600 pairs.

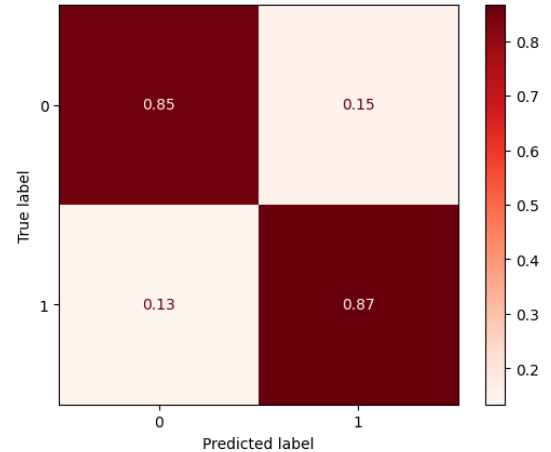
Fairness was a concern in this paper, especially considering the sensitive attributes of the graph. In order to explore this, the XGBoost model was adjusted using `fairlearn`'s threshold optimizer, which ultimately aimed to reduce the inequality through adjusting an acceptance threshold, respectively, for each attribute. The goal of this procedure is to enforce demographic parity and equalized odds.

## 1.8 Final results

**XGBoost Performance** Of the three models, XGBoost performed the best with a 5-fold cross-validation accuracy of  $0.822 \pm 0.12$  and  $0.873 \pm 0.01$  for training and testing, respectively. Figure 2(a) shows the ROC for this model on the test set, and table 7 shows the optimal parameters found. The



(a) ROC curve on the test set for the optimal XGBoost model. An AUC of 0.94 is achieved.



(b) Confusion matrix on the test set for the optimal XGBoost model. Values were normalized to represent respective rates.

**Fig. 2.** (a) Confusion matrix and (b) ROC curve on the test set for the optimal XGBoost model.

confusion matrix in figure 2(b) shows a breakdown of the prediction of the model. The XGBoost model appears to have a well balanced performance across, with similar TPR and TNR as well as FNR and FPR. This means that the classifier will predict positive and negative classes with similar rate. This balance was achieved by tuning the ration of positive to negative examples in training.

**Fairlearn Evaluation** An analysis of the resulting predictions from the highest performing XGBoost models was carried out to determine the level of fairness the model is capable of predicting. The selected criteria for overall fairness of the model was based on equalized odds, which denotes an equal true positive rate and equal false positive rate for each sensitive group, as well as the demographic parity, which denotes equal probability for selecting one outcome for one group against another. A fairlearn assessment was conducted in table ??, indicating an equalized odds difference of 0.54, as well as a demographic parity ratio of 0.12. This is an indication that equalized odds does not hold for the model selected.

**Table 1.** Fairlearn evaluation results for XGBoost model before and after threshold optimization.

Metric	Value	Metric	Value
Overall selection rate	0.57	Overall selection rate	0.30
Demographic parity difference	0.67	Demographic parity difference	0.24
Demographic parity ratio	0.12	Demographic parity ratio	0.69
False positive rate difference	0.5	False positive rate difference	0.22
False negative rate difference	0.55	False negative rate difference	0.32
Equalized odds difference	0.55	Equalized odds difference	0.32
Overall AUC	0.84	Overall AUC	0.62
AUC difference	0.26	AUC difference	0.18

**Table 2.** Before threshold optimizing

**Table 3.** After threshold optimizing

From tables 1 and 1, an improvement in the demographic parity ratio and equalized odds difference can be seen. Indeed, through the new adjusted model displayed an equalized odds difference of 0.32, an improvement on the original model's 0.55, and a demographic parity ratio of 0.69, as opposed to the original model's 0.12, indicating a more fair model. However, despite having a more fair model, the AUC score is significantly decreased, as the original AUC was 0.84, which dropped to an overall AUC of 0.62. Therefore, despite adjustments to the fairness of the model to accommodate the demographic parity and equalized odds, the overall performance of the model dropped.

## 1.9 Limitations

Although the XGBoost model performs well, it also comes with limitations. One of its main limitations is the risk of over-fitting [2]. Especially because of the relative small size of the training set, over-fitting is more likely. However, its tendency to over-fit can be attenuated by careful hyperparameter optimization; however given the scope of this paper, an extensive grid search was too computationally expensive.

Another pitfall of this model is its ensemble nature which leads to a black box type model. XGBoost does not allow for predictions to be explained like, for example, a decision tree. Given the nature of the link prediction problem and its relevant societal applications, this limitation can lead to ethical problems arising. Although, model agnostic method to explain model prediction exist to alleviate this concern. Tools like SHAP can be employed, but were not possible given the scope of this paper.

Some limitations also arise from the synthetic nature of the training data. Leakage was noticed in the node id columns, likely due to the method of graph generation. Indeed the node1 and node2 columns were noticed to have a significant features importance in the Random Forest, whereas this encoding should be, in theory, totally random. This can lead to a higher accuracy than would be possible given the same problem with natural data.

Additionally, from the fairness evaluation of the model, it was evident that the model was susceptible to some forms of bias. In this case, unequal odds for certain sensitive features. Aside from this, other concerns around privacy and transparency presented themselves. Given the black box nature of the model, there are some concerns about the transparency and explainability of the model, which were not evaluated, due to time and resources constraints. In addition, privacy issues around the sensitivity of the features being used for predictions. The concerns around the mentioned aspects of the model can further be mitigated in a larger scope of research.

## 2 Ethics essay

### 2.1 Introduction

In this paper, a link prediction algorithm was proposed, taking in an existing network, along with a list of sensitive features associated with nodes in the network. The proposed link prediction algorithm displayed strong results on its validation sets, a result of a variety of feature engineering techniques, model selection and fine-tuning. Predicting links in large networks has a wide variety of implications outside of the problem investigated in this paper. The solution is applicable in various fields and problems which may be represented as a network. Additionally, link prediction algorithms are ubiquitously used within networks of people specifically, given that within large social networks, features about each individual provide more information to models. Despite the strong performance of such link prediction algorithms and their effectiveness in deriving important information, the applications of link prediction algorithms may lend themselves to an array of ethical concerns. [6]

In this section of the paper, its social and ethical implications around fairness, privacy and transparency are critiqued. With regards to the model fairness, the given model is trained on a sparse dataset, assigning importance to each entry, based on certain features, and how well the features can explain the output. While this may yield accurate results, sensitive features may provide unfairness to marginalized groups for instance. Even in tweaking the model accordingly, with removing sensitive features, reducing thresholds for selection, or even selecting other models, there is the risk of losing performance, or even maintaining a similar level of discrimination through other factors. The mitigation of these risks proves difficult. This is due to the fact that models will continue to fit its training data, and even with the removal of certain biases, it is difficult to guarantee that other biases will not arise.

With regards to privacy, it is difficult to ensure that the model respects the aspect of privacy, especially in a social network space. In any case where the nodes within the network may represent individuals in a social web, The algorithm proposed may create issues around privacy. The model proposed achieved a validation accuracy of over 87%, only using information from the sensitive attributes of each node and the graph itself. Put into perspective, a simple model with the capabilities of predicting links based on only 6 features will almost certainly correctly predict whether a link exists between one node and another. In large social networks, where significantly more information exists around each individual, the implication of providing similar solutions is worrisome, in that with state of the art technology, algorithms will certainly be able to provide high confidence in determining links between individuals, which those individuals are not willing to make public. Furthermore, with the solution proposed above, high-performing models have the capacity to work both ways. In the example presented previously, the model is able to predict sensitive attributes of nodes, solely using the engineered features of each node, which used in a social network context, poses a severe threat to the privacy of those individuals.

The proposed solution also has significant implications with regards to transparency. The algorithm utilized was an XGBoost model, pulled from the SciKit-Learn package. The model was used as a black box, meaning there are various hyper parameters which could be set, however ultimately, the model is being utilized in such a way where inputs are fed through, and an output is produced. The black box nature of the model removes an aspect of transparency with what is being performed on the data. For instance, it is difficult to determine whether the model is providing less importance on certain features deliberately, or if the model is creating discrimination in its decision power. Therefore if the model scales up, or is used in highly decisive settings, there may arise the issue of its lack of transparency.

### 2.2 Elaboration

Concerns over fairness are festering in the current technological climate where over 5.17 billion people worldwide actively participate in social media. A large amount of information is collected from each of these individuals, and in a capitalistic environment, in which success is measured in accumulating high rates of retention and acquisition, the applications of these algorithms is concerning in how those who possess collected data have power in using such algorithms to their own benefit, even if the model being utilized is violating regulations of fairness and discrimination.

The proposed solution in this paper performs strongly through the processing of data, engineering features, and tuning models. With a training set of only approximately 13,000 entries, and beginning with only 6 different features, the accumulated accuracy of the proposed model is sufficient in predicting links. Put into perspective, in a setting where 5.17 billion willing individuals make their data public, the training data for such models, as well as the budget for state-of-the-art models, will result in models which can infer a staggering amount of detailed personal insights.

The effects of biased link prediction algorithms are highly detrimental in the longitudinal time frame since they might lead to the creation of class structures and entrenching of societal stereotypes. As the nature of algorithms also determines the kind of social ties made by persons, then they can promote already existing societal biases which may kick-start a cycle of discrimination that is hard to break.

One notable instance in particular is centered around a lawsuit for the popular social media outlet, TikTok. The platform has been accused of leveraging its "addictive algorithm" as part of a profit-driven model that prioritizes user retention, particularly targeting vulnerable groups such as children.[1] The lawsuit, filed by 13 states and the District of Columbia, claims that TikTok falsely advertises its content moderation and safety policies for minors while intentionally encouraging engagement through its algorithm. The TikTok recommendation algorithm is catered to its users, and more specifically through their engagement with certain posts, and other interactions. With the amount of users actively engaging with the application, the data collection practices include the collection of information on children. The sentiment of collecting data on the social media behavior has both its benefits and its consequences. In an ideal setting, understanding the behavior of children through their interactions would allow the application to monitor content, and tailor the content such that explicit, or harmful content is suppressed. However in a setting where the primary objective of the application is to maximise interaction, this may not be the case. Perhaps the intention of the recommender algorithm is not solely to harm vulnerable groups. However if the algorithm aims to encourage retention, and maximise profit, it may undeliberately target these vulnerable groups.[3] As a result, consequences to vulnerable groups, in this case, minors and children using the app, are severe. It is shown that minors who spend extensive time on the social media application display signs of mental health decline due to the unrealistic expectations of the content being shown in the platform.

Furthermore, through the exploitation of such an algorithm, the risk of introducing influential or political content to children is harmful in their mental development. While the intention may not be solely to harm vulnerable groups, the algorithm's primary objective to maximize interaction and profit could lead to unintended consequences for these populations. Minors using the app have shown signs of mental health decline due to unrealistic expectations set by the content displayed. Furthermore, the risk of exposing children to influential or political content poses additional concerns regarding their mental development.

In the case of this paper, modifying the model such that it would reduce the risk of unequal odds came at the cost of model accuracy. For a company such as TikTok, the trade-off between having a more fair model with the potential loss of accuracy could mean a sacrifice in their financial gains. With the sheer computing power available, as well as the resources TikTok has at its disposal, it likely possesses the capabilities to mitigate these risks. Nonetheless, actions taken to improve fairness of its algorithm seem to be non-existent, as long as its primary objective is to maximise interaction, retention and profit. To mitigate these risks, developers of link prediction algorithms should consider implementing fairness-aware models and establishing ongoing monitoring and user feedback mechanisms, fostering a more responsible approach to AI development.

## 2.3 Mitigation

The mitigation of such risks around fairness proves to be a difficult task, where solutions often differ on a case-to-case basis. Model behavior is difficult to control, especially if the models gain more accuracy in exploiting correlations with sensitive features in its data.

Tatineni [6] proposes various mitigation techniques, considering the roles and responsibilities of each involved entity, namely, the developers, the government, and the organizations. Tatineni suggests ethical impact assessments implemented in the workflow of product development, ensuring that model fairness is evaluated frequently, ensuring it does not create any bias when the model is deployed. Secondly, they suggest agile ethical frameworks, which account for adaptive ethical frameworks, which can easily evolve alongside technological and social advancements. This is vital to large projects, in adapting how they view fairness overall, as new findings are uncovered. Thirdly, they encourage wide-scale collaboration, employing experts, researchers, policymakers, in consistently challenging the ethical aspects of the model, creating new goals and objectives to strive towards.

There are a wide variety of methods used in the field of data science where biases towards sensitive features are mitigated during the pre-processing, in-processing and post-processing of data being used in the model. The process of conducting these mitigation techniques is non-trivial, as biases around different problems arise in different ways. If biases already exist in the training data, the model will then fit to these biases. However more complicated models for instance may even infer these underlying biases,

even if it is not directly present in the training set. In preparing and pre-processing of the data, some methods for eliminating bias include data cleaning, resampling, or the introduction of synthetic features. In many cases, this is done in order to alter the training set such that decisions made by the model are less influenced by sensitive features. One example is in predicting the probability of defaulting on a credit card payment, based on features of each individual user. If the model is more likely to predict a default payment for individuals of a certain gender, race or marital status for instance, one can synthesize features such that the probability of default is equal amongst individuals from each sensitive group. This allows models or algorithms to ultimately predict a default in a more fair manner, reducing the risk of predicting more highly for certain groups.

In reference to the TikTok case, there are a variety of methods they may take in order to tackle some of these risks. For instance, they may apply data anonymization or data minimization techniques to ensure that sensitive attributes such as age are either removed, masked, or down-weighted before being processed by the algorithm. For this issue in particular, this may remove bias from the algorithm in targeting younger users. Additionally, various post-processing techniques exist in mitigating these biases. In the case of TikTok, this may include the analysis of the final output of the algorithm, and the adjustment of the content to groups which are knowingly being targeted. Furthermore, one may also consider policy-based interventions. This may include the implementation of further regulations or privileges on the application, restricting the time spent, perhaps through parental control and age-verification. This in turn, may limit the amount of data collected on these vulnerable groups mitigating both fairness and privacy risks.

On a larger scale, there also exist a variety of government regulations to be implemented to mitigate these risks [6], and hold TikTok accountable for the issues caused by their application, ultimately incentivizing them to adjust their algorithm to comply with more fair practices. The example of TikTok, and the methods proposed in tackling the issue of fairness on their platform, are not only applicable to the example itself, but rather in problems across all fields. In the field of data science, models are rewarded according to an objective, and how well those models are equipped in achieving these results. Fundamentally, if an entity aims to maximise its profits, its model will achieve it by any means necessary, even at the cost of the well-being of vulnerable groups. In a general sense, it is of utmost importance to make the shift from solely profit-based models, and understand the importance of fairness in today's capitalistic climate, even if it comes at the cost of performance.

Evidently, with the example presented, it may be observed that the issues which may arise from irresponsible use of data, can be severe, and the prevention of these issues may come at the trade-off of some performance.

## 2.4 Conclusion

The significance of issues which may arise from failing to comply with fairness procedures is profound, as it can lead to severe consequences, which were elaborated on using the TikTok case. The magnitude of such consequences may severely impact individuals in a negative manner. Overall, the research conducted in this study was aimed at critiquing instances where fairness, privacy, and transparency are violated when implementing large-scale models similar to the works presented in the technical section of the research paper. Along with the exponential increase in capacity of machine learning techniques and models, it is becoming increasingly more trivial in achieving high model accuracies. Therefore, it is of utmost importance to grasp and understand the ethical implications of the access to such models, and technologies, in order to minimize the harm caused to individuals. In this paper, a variety of mitigation techniques, both technology-based, and policy-based, were presented. In this paper, a number of mitigating measures, both of a technological nature and policy escalation, were offered. These include approaches like data cleanup, data resampling, and the deployment of novel synthetic features for bias removal, as well as legal and institutional measures which prohibit and punish abuse. The application of these four strategies is critical in mitigating the ethical concerns on the use of machine learning applications so that technology can be safely and positively used in society. Fairness, privacy and transparency are the aspects that stakeholders must consider to ensure trust in AI systems and equality of opportunities for all.

Ultimately, these mitigation strategies remind us of the crucial roles and social responsibilities that data scientists possess in ensuring that technological advancements do not come at the expense of vulnerable groups, highlighting the need for a shift from profit-driven models to those that value fairness and ethical considerations.



### 3 Appendix A

#### 3.1 Hyperparameters Grid

Hyperparameter	Values
n_estimators	{10, 20, 50, 100, 200}
max_depth	{None, 5, 10}
min_samples_split	{2, 5, 10}
min_samples_leaf	{1, 2, 4}

**Table 4.** Hyperparameters for Random Forest

Hyperparameter	Values
lr__C	{0.001, 0.01, 0.1, 1, 10, 100, 1000}

**Table 5.** Hyperparameters for Logistic Regression Model

Hyperparameter	Values
n_estimators	{50, 100, 200}
max_depth	{3, 4, 5}
learning_rate	{1, 0.5, 0.2, 0.1}
subsample	{0.5, 0.6, 0.7}
gamma	{0, 0.01, 0.1}
reg_alpha	{0, 0.001}
reg_lambda	{0.5, 1, 1.5}

**Table 6.** Hyperparameters for XGBoost Model

Hyperparameter	Optimal Value
gamma	0.01
learning_rate	0.1
max_depth	3
n_estimators	100
reg_alpha	0
reg_lambda	1.5
subsample	0.6

**Table 7.** Optimal Hyperparameters for XGBoost Model

## References

- [1] Chase DiBenedetto. *TikTok sued by more than a dozen states for allegedly addictive algorithm*. Oct. 2024. URL: [https://mashable.com/article/dozen-states-suing-tiktok-harming-young-people?test\\_uuid=01iI2GpryXngy77uIpA3Y4B&test\\_variant=a](https://mashable.com/article/dozen-states-suing-tiktok-harming-young-people?test_uuid=01iI2GpryXngy77uIpA3Y4B&test_variant=a).
- [2] Christina Ellis. *XGBoost overfitting*. Crunching the Data, Aug. 2022. URL: <https://crunchingthedata.com/xgboost-overfitting/>.
- [3] Michal Lavi et al. *TARGETING CHILDREN: LIABILITY FOR ALGORITHMIC RECOMMENDATIONS*. URL: [https://aulawreview.org/wp-content/uploads/2024/09/Lavi.to\\_.Printer.pdf](https://aulawreview.org/wp-content/uploads/2024/09/Lavi.to_.Printer.pdf).
- [4] Y V Nandini, T. Jaya Lakshmi, and Murali Krishna Enduri. “Link Prediction in Complex Networks: An Empirical Review”. In: *Smart innovation, systems and technologies* (Jan. 2023), pp. 57–67. DOI: 10.1007/978-981-99-6706-3\_5. (Visited on 10/13/2024).
- [5] Zhan Su et al. “Link prediction in recommender systems based on vector similarity”. In: *Physica A: Statistical Mechanics and its Applications* 560 (Dec. 2020), p. 125154. DOI: 10.1016/j.physa.2020.125154. (Visited on 04/08/2021).
- [6] Sumanth Tatineni. “Ethical Considerations in AI and Data Science: Bias, Fairness, and Accountability”. In: 10.1 (2019), pp. 11–20. URL: [https://www.researchgate.net/profile/Sumanth-Tatineni/publication/377701616\\_Ethical\\_Considerations\\_in\\_AI\\_and\\_Data\\_Science\\_Bias\\_Fairness\\_and\\_Accountability/links/65b35fff34bbff5ba7c4d023/Ethical-Considerations-in-AI-and-Data-Science-Bias-Fairness-and-Accountability.pdf](https://www.researchgate.net/profile/Sumanth-Tatineni/publication/377701616_Ethical_Considerations_in_AI_and_Data_Science_Bias_Fairness_and_Accountability/links/65b35fff34bbff5ba7c4d023/Ethical-Considerations-in-AI-and-Data-Science-Bias-Fairness-and-Accountability.pdf).
- [7] Haixia Wu et al. “Link Prediction on Complex Networks: An Experimental Survey”. In: *Data Science and Engineering* 7 (June 2022), pp. 253–278. DOI: 10.1007/s41019-022-00188-2.