# Fairwalk: Towards Fair Graph Embedding

Tahleen Rahman, Bartlomiej Surma, Michael Backes and Yang Zhang
@ CISPA Helmholtz Center for Information Security

---

D6: Online Social Networks and Media; paper presentation

Papazis Stergios (483)

Zisopoulos Georgios (505)

Wednesday, 17 January 2024

Department of Computer Science & Engineering
School of Engineering
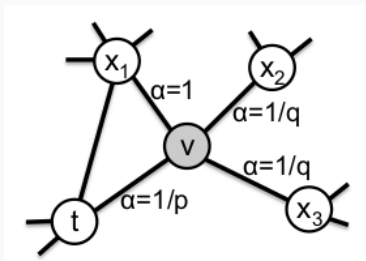University of Ioannina

## Overview

- Assume a graph task that requires ML
  (e.g. building a link/friend recommendation system)
- ML algorithms can only be used with vector representation
  ⇒ graph embedding are needed
- Node2vec generates embeddings, *but*:
    - fairness through unawareness
    - ignores minority structure
    - creates echo chambers, bubble effect
- How Node2vec is unfair and how to modify to fix it
  (Fairwalk)

# Nodes2vec

- state of the art algorithm
- captures structure of network (node similarity, hubs, neighbourhood structure)
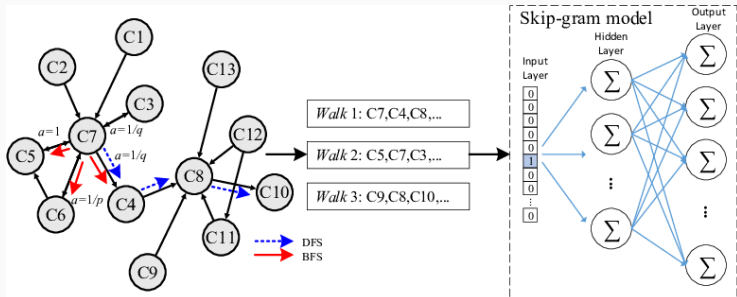- extendable to edge embedding

Walk generation

- second order random walks
- p: search width
- q: search depth
- parameters: d, walk_num, walk_len

# Nodes2vec

Train a ML model that creates the vector representations:

1. sample node neighbourhoods through random walks;
2. maximise likelihood that a node neighbourhood appear given the feature vector of the node

$$\max_f \sum_{u \in V} \log P\left(N_S(u) \mid f(u)\right)$$

- How to measure fairness in this setting?

## Fairness

- How to measure fairness in this setting?
- **Sensitive attribute**: $S$   (e.g. g: gender, r: race, ... )

## Fairness

- How to measure fairness in this setting?
- **Sensitive attribute**: $S$ (e.g. g: gender, r: race, ... )
- **Group** $G_{ij}^S$ : pair of users with attribute values $i$, $j$

## Fairness

- How to measure fairness in this setting?
- **Sensitive attribute**: $S$    (e.g. g: gender, r: race, … )
- **Group** $G_{ij}^S$ : pair of users with attribute values $i, j$
- **Number of recommendations for** $G_{ij}^S$ :

$$N_{rec}(G_{ij}^S) := \left| \left\{ (u, v) \in G_{ij}^S : v \in \rho(u) \right\} \right|$$

## Fairness

- How to measure fairness in this setting?
- **Sensitive attribute**: $S$  (e.g. g: gender, r: race, … )
- **Group** $G_{ij}^S$ : pair of users with attribute values $i, j$
- **Number of recommendations for** $G_{ij}^S$ :

$$N_{rec}(G_{ij}^S) := \left| \left\{ (u, v) \in G_{ij}^S : \ v \in \rho(u) \right\} \right|$$

- **Acceptance rate**:

$$\mathsf{P}(G_{ij}^S) = \frac{N_{rec}(G_{ij}^S)}{\left| G_{ij}^S \right|}$$

fraction of $j$-user suggested to $i$-users over total number of users with $i, j$ attributes

- **Statistical Parity** or **Demographic Parity** or **Independence**: acceptance rates of the candidates from all groups should be somewhat equal
- bias wrt $G_{ab}^S$, $G_{cd}^S$ :

$$\text{bias}^{SI}(G^S) := P(G_{ab}^S) - P(G_{cd}^S)$$

- bias for multiple groups :

$$\text{bias}^{SI}(G^S) := \text{Var}\left(\left\{P(G_{ij}^S) : \ G_{ij}^S \in G^S\right\}\right)$$

Equality of Representation

- **network level**:
$$\text{bias}^{\text{ERg}}(G^S) := \text{Var}\left(\left\{ N_{rec}(G_{ij}^S) : \ G_{ij}^S \in G^S \right\}\right)$$

- **user level**:
$$\text{bias}^{\text{ERu}}(z) := \frac{1}{|Z^S|} - \frac{1}{|U|} \sum_{u \in U} z\text{-share}(u)$$

  - measures recommendation fairness *independent* of the ground truth (existing friendships contain bias)
  - fair fraction minus average z-share over all users
  - positive: underrepresented, negative: overrepresented

z-share: fraction of recommended users with attribute z
$$z\text{-share}(u) := \frac{|\rho_z(u)|}{|\rho(u)|}$$

## Node2vec unfairness example

Dataset

- nodes: *Instagram users* in London & LA
- edges: *mutual follows*
- sensitive attributes: *gender*, *race*
  (autogenerated from profile photos)
- is biased

Groups

- $G_{00}^{g}$, $G_{01}^{g}$, $G_{10}^{g}$, $G_{11}^{g}$,
  (0: female, 1: male)
- $G_{00}^{r}$, $G_{01}^{r}$, $G_{10}^{r}$, $G_{11}^{r}$, $G_{12}^{r}$, $G_{21}^{r}$, $G_{22}^{r}$, $G_{02}^{r}$, $G_{20}^{r}$
  (0: african, 1: caucasian, 2: asian)

# Node2vec unfairness example

|                  | LA      | London  |
|------------------|---------|---------|
| No. users        | 82,607  | 53,902  |
| No. social links | 482,305 | 165,184 |
| gender 0         | 62.6%   | 62.3%   |
| gender 1         | 37.4%   | 37.7%   |
| race 0           | 21.9%   | 15.9%   |
| race 1           | 72.2%   | 80.7%   |
| race 2           | 5.9%    | 3.4%    |

Table 2: Statistics of both datasets.

|                      | Gender groups | | | | Race groups | | | | | | | | |
|----------------------|-------|-------|-------|-------|------|-------|-------|-------|------|------|------|------|------|
| $i - j$ for $G_{ij}$ | 0-0   | 0-1   | 1-0   | 1-1   | 0-0  | 0-1   | 1-0   | 1-1   | 0-2  | 1-2  | 2-0  | 2-1  | 2-2  |
| LA                   | 37.78 | 21.00 | 21.99 | 19.21 | 5.97 | 12.55 | 12.65 | 57.92 | 1.17 | 3.87 | 1.16 | 3.74 | 0.96 |
| London               | 38.96 | 18.19 | 20.74 | 22.10 | 3.55 | 9.52  | 9.83  | 70.31 | 0.47 | 2.57 | 0.56 | 2.76 | 0.41 |

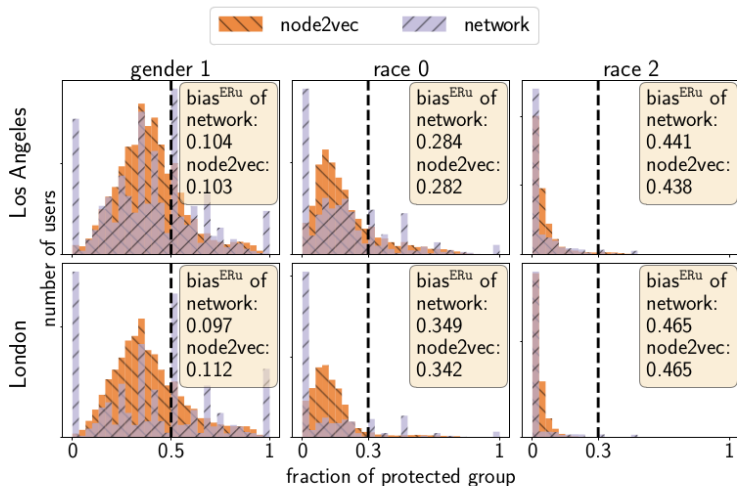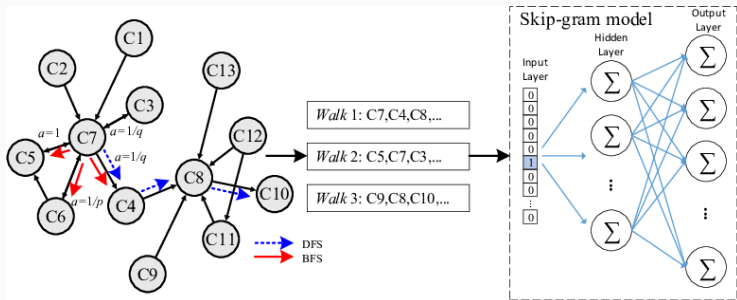Table 1: Percentage of existing friendships in each group in our original dataset

Figure 1: $z$-share distributions of node2vec and original network. The vertical line shows the *fair fraction* (0.5 and 0.3)

Modify random walk generation

# Fair walk generation

**Algorithm 1** Fair random walk trace generation

1: **procedure** RAND_WALK($U, \omega, \texttt{walk\_num}, \texttt{walk\_len}$)
2:      traces $\leftarrow$ empty_list
3:      **for all** $u \in U$ **do**
4:          **for** $i \leftarrow 0, \texttt{walk\_num}$ **do**
5:              trace $\leftarrow$ empty_list
6:              $u_1 \leftarrow u$
7:              **for** $j \leftarrow 0, \texttt{walk\_len}$ **do**
8:                  trace.append($u_1$)
9:                  $\mathcal{Z}_u \leftarrow \{z : z \in \mathcal{Z} \wedge |\omega_z(u_1)| > 0\}$
10:                 $z_1 \xleftarrow{R} \mathcal{Z}_u$
11:                 $v \xleftarrow{R} \omega_{z_1}(u_1)$
12:                 $u_1 \leftarrow v$
13:              **end for**
14:              traces.append(trace)
15:          **end for**
16:      **end for**
17:      **return** traces
18: **end procedure**

- Node2vec mirrors underlying distribution, not enough info about minorities
- Fair walks create higher network diversity, greater representation of minorities
  $\Rightarrow$ ML gets better understanding of sensitive attributes


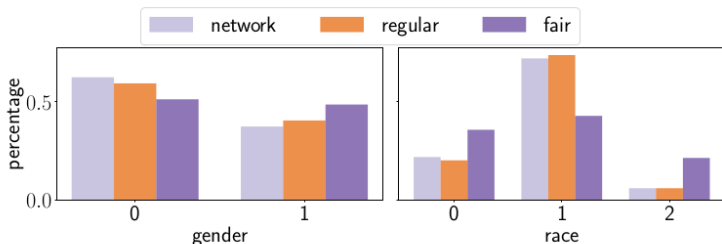
Figure 2: Ratio of each gender and race in the original network and regular and fair random walk traces in Los Angeles dataset

## Experiment

Experimental Setup

- 5 iterations
- 80% training set
- 20% test set
- generate node embeddings
- train random forest on embedded pairs of users in order to learn presence or absence of friendships (links)
- recommend $k$% most similar non-friend users

## Dataset reminder

Dataset

- nodes: *Instagram users* in London & LA
- edges: *mutual follows*
- sensitive attributes: *gender*, *race*
  (autogenerated from profile photos)
- is biased

Groups

- $G_{00}^{g}$, $G_{01}^{g}$, $G_{10}^{g}$, $G_{11}^{g}$,
  (0: female, 1: male)
- $G_{00}^{r}$, $G_{01}^{r}$, $G_{10}^{r}$, $G_{11}^{r}$, $G_{12}^{r}$, $G_{21}^{r}$, $G_{22}^{r}$, $G_{02}^{r}$, $G_{20}^{r}$
  (0: african, 1: caucasian, 2: asian)

|                  | LA      | London  |
|------------------|---------|---------|
| No. users        | 82,607  | 53,902  |
| No. social links | 482,305 | 165,184 |
| gender 0         | 62.6%   | 62.3%   |
| gender 1         | 37.4%   | 37.7%   |
| race 0           | 21.9%   | 15.9%   |
| race 1           | 72.2%   | 80.7%   |
| race 2           | 5.9%    | 3.4%    |

Table 2: Statistics of both datasets.

| $i - j$ for $G_{ij}$ | Gender groups | | | | Race groups | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 0-0 | 0-1 | 1-0 | 1-1 | 0-0 | 0-1 | 1-0 | 1-1 | 0-2 | 1-2 | 2-0 | 2-1 | 2-2 |
| LA | 37.78 | 21.00 | 21.99 | 19.21 | 5.97 | 12.55 | 12.65 | 57.92 | 1.17 | 3.87 | 1.16 | 3.74 | 0.96 |
| London | 38.96 | 18.19 | 20.74 | 22.10 | 3.55 | 9.52 | 9.83 | 70.31 | 0.47 | 2.57 | 0.56 | 2.76 | 0.41 |

Table 1: Percentage of existing friendships in each group in our original dataset

| | | LA | | London | |
|---|---|---|---|---|---|
| | | gender | race | gender | race |
| ERg | regular | $1.3e^{10}$ | $2.5e^7$ | $6.5e^9$ | $2.4e^7$ |
| | fair | $0.8e^{10}$ | $1.9e^7$ | $4.8e^9$ | $1.9e^7$ |
| SI | regular | $4.7e^{-9}$ | $1.4e^{-12}$ | $1.1e^{-8}$ | $7.1e^{-11}$ |
| | fair | $1.7e^{-9}$ | $0.4e^{-12}$ | $0.2e^{-8}$ | $2.8e^{-11}$ |

Table 3: bias$^{\text{SI}}$ and bias$^{\text{ERg}}$ for both cities (lower, the better)

Figure 3: Fraction of recommended users pairs out of all possible pairs in each group. The x-axes marks the Type-2 groups $G_{z_i, z_j}$ with the corresponding $i - j$
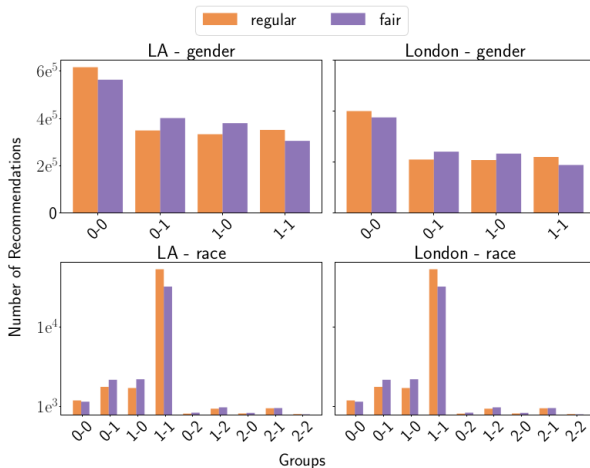
Figure 4: Number of recommended users pairs from each group. The x-axes marks groups $G_{ij}$ with the corresponding $i - j$

|  |  | gender | | race | | |
|---|---|---|---|---|---|---|
|  |  | 0 | 1 | 0 | 1 | 2 |
| LA | network | 0.104 | 0.104 | 0.117 | 0.392 | 0.275 |
|  | node2vec | 0.103 | 0.103 | 0.115 | 0.387 | 0.272 |
|  | fairwalk | 0.068 | 0.068 | 0.054 | 0.288 | 0.234 |
| London | network | 0.097 | 0.097 | 0.183 | 0.481 | 0.298 |
|  | node2vec | 0.112 | 0.112 | 0.176 | 0.474 | 0.298 |
|  | fairwalk | 0.095 | 0.095 | 0.135 | 0.417 | 0.282 |

Table 4: Bias by *Equality of Representation* at user level for both genders and all three races (lower, the better).
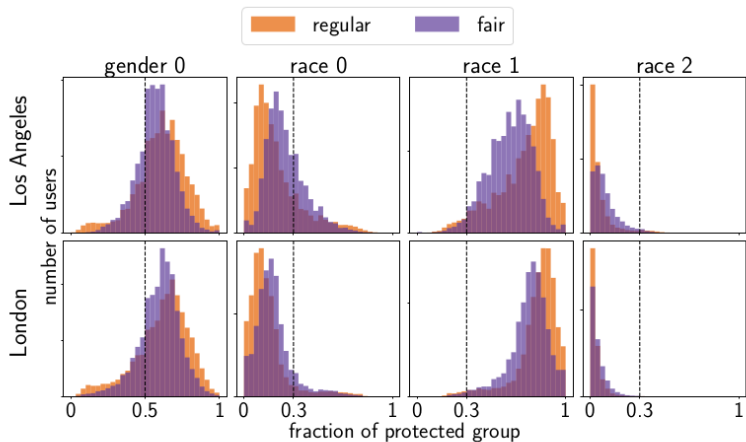
Figure 5: $z$-share distributions of node2vec and *Fairwalk*. The vertical line shows the *fair fraction*.
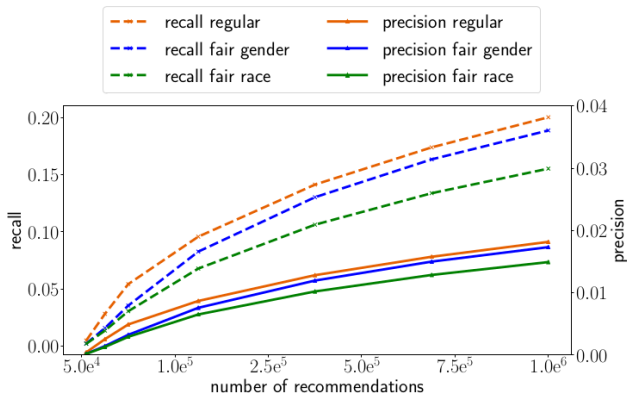
Figure 6: Precision and recall for different number of recommendations.

$$\text{recall} = \frac{TP}{TP + FN}, \qquad \text{precision} = \frac{TP}{TP + FP}$$

Thanks for your attention!

# References

[1]    Aditya Grover and Jure Leskovec. "Node2vec: Scalable Feature Learning for Networks". In: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*. KDD '16. San Francisco, California, USA: Association for Computing Machinery, 2016, pp. 855–864. DOI: 10.1145/2939672.2939754.

[2]    Tahleen Rahman et al. "Fairwalk: Towards Fair Graph Embedding". In: *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*. International Joint Conferences on Artificial Intelligence Organization, July 2019, pp. 3289–3295. DOI: 10.24963/ijcai.2019/456.