

Outline for today

- ▶ What is fairness?
- ▶ Fair Link Prediction
- ▶ Fair Document Classification
- ▶ Concluding Remarks

What does it mean to be fair ?

It depends ...

A matter of perspective

Social Sciences

A system is fair if people perceive it as being fair



A matter of perspective

Philosophy

Fairness asks what is morally right, focusing on how decisions should be made.



A matter of perspective

Law

Fairness is about preventing discrimination and limiting *unjustified* unequal outcomes.



A matter of perspective

Quantitative Fields

Fairness is addressed through a pure quantitative perspective



Fairness in Supervised Learning

($\underbrace{\text{features}}_X, \underbrace{\text{sensitive attribute}}_A, \underbrace{\text{label}}_Y$) P on $X \times A \times Y$

Goal: learn a function $h : X \rightarrow Y$ such that

$$\min E_{(X,Y) \sim P} [h(X), Y] \quad \text{s.t.} \quad \text{Fairness}(h) \leq \epsilon,$$

where $\text{Fairness}(h)$ can be,

- ▶ $DP : |P(h(X) = 1 / A = 0) - P(h(X) = 1 / A = 1)|$
- ▶ $EO : |P(h(X) = 1 / A = 0, Y = 1) - P(h(X) = 1 / A = 1, Y = 1)|$

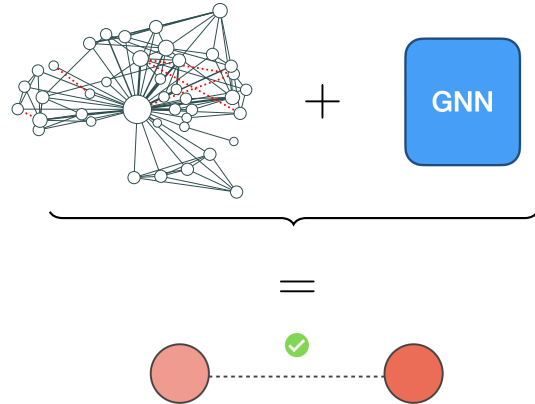
NLP

x_i : she is able to assess, diagnose [. . .]
 y_i : nurse
 a_i : women



Graph

x_i : pair of nodes (u, v)
 y_i : existence of a link
 a_i : ?



What could go wrong?

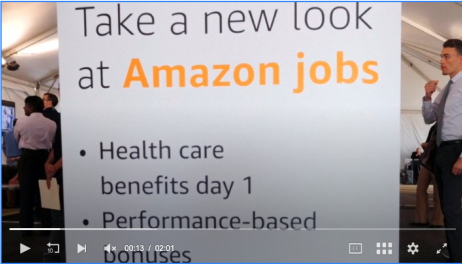
What could go wrong?

NLP

Insight - Amazon scraps secret AI recruiting tool that showed bias against women

By Jeffrey Dastin

October 11, 2018 2:50 AM GMT+2 · Updated October 11, 2018



Graph

Outline for today

- ▶ What is fairness?
- ▶ **Fair Link Prediction**
- ▶ Fair Document Classification
- ▶ Concluding Remarks

Fair Link Prediction¹

$g = (V, E, F, A)$, where $V = \{v_i\}_{i=1}^m$, $E \subseteq V \times V$, $\mathbf{F} \in \mathbb{R}^{m \times d}$ and $A \in \mathbb{R}^m$

Goal: learn a **fair** link prediction model $h : V \times V \rightarrow [0, 1]$.

☀ Differentiate intra vs. cross group edges.

¹Laclau C. et al. All of the Fairness for Edge Prediction with Optimal Transport. AISTATS 2021

Fair Link Prediction¹

$g = (V, E, F, A)$, where $V = \{v_i\}_{i=1}^m$, $E \subseteq V \times V$, $\mathbf{F} \in \mathbb{R}^{m \times d}$ and $A \in \mathbb{R}^m$

Goal: learn a **fair** link prediction model $h : V \times V \rightarrow [0, 1]$.

☀ Differentiate intra vs. cross group edges.

Dyadic Fairness (DP)

We can define DP for link prediction as

$$DP : \mathbb{P}(h(V, V) = 1/A \mid A = 1) - \mathbb{P}(h(V, V) = 1/A \mid A = 0)$$

¹Laclau C. et al. All of the Fairness for Edge Prediction with Optimal Transport. AISTATS 2021

Connection between homophily and fairness ³

The homophily mechanism

People generally prefer to befriend others similar to themselves. ²

²[Byrne 1971; McPherson et al. 2001]

³Marey L., Perez M., Viard T., Laclau C., TopoFair: Linking Topological Bias to Fairness in Link Prediction Benchmarks, arxiv

Connection between homophily and fairness ³

The homophily mechanism

People generally prefer to befriend others similar to themselves. ²

The homophily mechanism is at the origin of segregation in online networks.

²[Byrne 1971; McPherson et al. 2001]

³Marey L., Perez M., Viard T., Laclau C., TopoFair: Linking Topological Bias to Fairness in Link Prediction Benchmarks, arxiv

Revisiting Dyadic Fairness⁴

- | $A = f$ blue, red g
- | We consider two edge predictors ($h^{(1)}$ and $h^{(2)}$), each predicting two edges.
- | All graphs have the same assortativity (homophily proxy).
- | Both predictors have the same disparate impact

⁴Marey L., Viard T., Laclau C. k-hop Fairness: Addressing Disparities in Graph Link Prediction Beyond First-Order Neighborhoods, under review.

Defining k-hop Exposure

For $v \in V$ and $k \in \mathbb{N}$, we define the k-hop neighborhood of v as

$$N^{(k)}(v) = \{v' \in V \mid (v, v') = k\};$$

k-hop node exposure

For $s \in \{0, 1, g\}$, $k \in \mathbb{N}$, and $v \in V$ and $\hat{Y} = h(V; V^0)$ we define

$$f_s^{(k)}(v) = \begin{cases} \sum_{v' \in N^{(k)}(v)} \hat{Y}_{f_{S^0=sg}} & \text{if } N^{(k)}(v) \neq \emptyset; \\ 0 & \text{otherwise;} \end{cases}$$

$$f_b^1(v) = 5=11; \quad f_r^1(v) = 6=11$$

$$f_b^2(v) = 5=11; \quad f_r^2(v) = 5=11$$

Connection with DP

We can rewrite DP across k (nothing fancy just the law of total probability).

$$DP = \sum_{k>0} \sum_{v \in V} !^{(k)}(v) \frac{f_{\text{same}}^{(k)}(v)}{P(S = S^0)} \frac{f_{\text{di}}^{(k)}(v)}{P(S \neq S^0)}$$

with

- | $f_{\text{same}}^{(k)}(v)$ and $f_{\text{di}}^{(k)}(v)$ the within and cross group exposure;
- | $!^{(k)}(v) = P(V = v; (V; V^0) = k)$

Connection with DP

We can rewrite DP across k (nothing fancy just the law of total probability).

$$DP = \sum_{k>0} \sum_{v \in V} !^{(k)}(v) \frac{f_{\text{same}}^{(k)}(v)}{P(S = S^0)} \frac{f_{\text{di}}^{(k)}(v)}{P(S \neq S^0)}$$

with

- | $f_{\text{same}}^{(k)}(v)$ and $f_{\text{di}}^{(k)}(v)$ the within and cross group exposure;
- | $!^{(k)}(v) = P(V = v; (V; V^0) = k)$

K-hop Fairness

k-hop fairness.

For a fixed graph distance $k \geq 1$, the fairness gap of a link predictor h is defined as

$$NF^{(k)}(h) = \max_{s_2 \in S} \max_{s_1, s_2 \in S} \frac{\binom{k}{s_1} (h)}{\binom{k}{s_2} (h)} :$$

In the binary case, $NF^{(k)}(h) = \max_{s \in \{0,1\}} \frac{\binom{k}{s} (h)}{\binom{k}{s_1} (h)} :$

1. Look at which group is the most unequally treated.
2. For that group, measure where the access imbalance is the largest.

K-hop Fairness

k-hop fairness.

For a fixed graph distance $k \geq 1$, the fairness gap of a link predictor h is defined as

$$NF^{(k)}(h) = \max_{s_2 \in S} \max_{s_1, s_2 \in S} \frac{P_{s_1, s_2}^{(k)}(h)}{P_{s_2, s_1}^{(k)}(h)} :$$

In the binary case, $NF^{(k)}(h) = \max_{s_2 \in \{0,1\}} \frac{P_{s_2, 0}^{(k)}(h)}{P_{s_2, 1}^{(k)}(h)} :$

1. Look at which group is the most unequally treated.
2. For that group, measure where the access imbalance is the largest.

Two strategies

rewiring the adjacency matrix or altering the edge probability.

Our findings

Next problem: Choice vs Induced Homophily

Next problem: Choice vs Induced Homophily

$$\text{Homophily} = \underbrace{\text{Choice}}_{\substack{\text{intrinsic preferences} \\ \text{Who I want to connect to}}} + \underbrace{\text{Induced}}_{\substack{\text{exposure \& structural constraints} \\ \text{Who I am exposed to}}}$$

Dynamic implication

Link prediction shapes the graph and modifies exposure patterns.

Choice vs Induced Homophily a Dynamic Process⁵

☀ Model group-pair interactions as a multidimensional Hawkes process.

We express Pinteraction (i; j) in [t; t + dt) as

$$ij(t) = \underbrace{ij}_{\text{choice}} + \int_0^t \underbrace{z_{(k;l)}(i;j)}_{\text{induced}} N_{kl}(ds)$$

where $z_{(k;l)}(i;j)$ (u) can encode the score of a link prediction model.

⁵Perez M., Romero R., Lijjt J., Laclau C., How Predicted Links Influence Network Evolution: Disentangling Choice and Algorithmic Feedback in Dynamic Graphs, under review.

Choice vs Induced Homophily a Dynamic Process⁵

☀ Model group-pair interactions as a multidimensional Hawkes process.

We express Pinteraction (i; j) in [t; t + dt) as

$$ij(t) = \underbrace{ij}_{\text{choice}} + \int_0^t \underbrace{z_{(k;l)|(i;j)}(t-s)}_{\text{induced}} N_{kl}(ds)$$

where $z_{(k;l)|(i;j)}(u)$ can encode the score of a link prediction model.

Given $w(t)$ and $c(t)$ the aggregated within-group and cross-group intensities:

$$B_{\text{inst}}(t) = \frac{w(t)}{w(t) + c(t)}; \quad B_{\text{inst}}(t) \in [0; 1]$$

⁵Perez M., Romero R., Lijjt J., Laclau C., How Predicted Links Influence Network Evolution: Disentangling Choice and Algorithmic Feedback in Dynamic Graphs, under review.

Long-Term Behavior of Instantaneous Bias

How does $B_{\text{inst}}(t)$ evolve over time?

From stochastic dynamics to average evolution

Taking expectations in the Hawkes system yields:

$$B_{\text{inst}}(t) = B_{\text{inst}} + \sum_{(k;l)} \int_0^t B_{\text{inst};kl}(t-s) \kappa_{kl}(s) ds$$

- | The micro-dynamics are stochastic.
- | The mean-field system is deterministic.
- | Its stability governs the convergence of $B_{\text{inst}}(t)$.

An example: professional network with two groups

An example: professional network with two groups

Main Theoretical Insights

- | **Stability depends on reinforcement strength.**
If excitation is weaker than decay, interaction intensities converge to a stable regime.
- | **Strong reinforcement amplifies asymmetries.**
When excitation dominates, small initial imbalances can grow and persist.
- | **Speed of convergence matters.**
If the system converges slowly, transient biases may dominate observed behavior.

Outline for today

- | What is fairness?
- | Fair Link Prediction
- | Fair Document Classification
- | Concluding Remarks

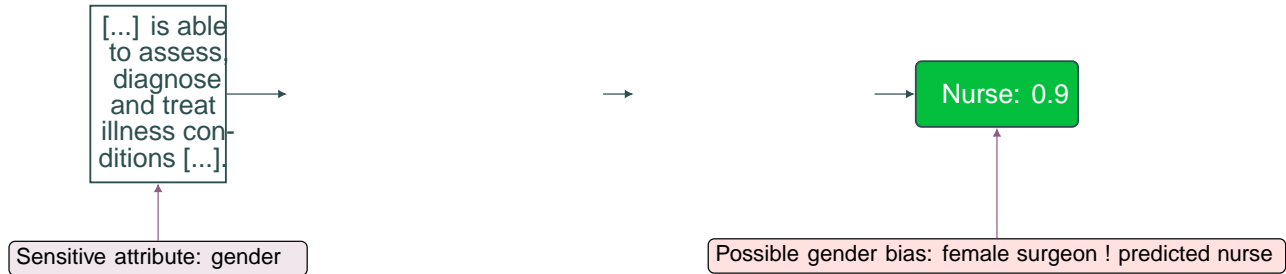
Bias lives in pre-trained models

- | Word embeddings trained on large corpora encode social stereotypes
- | Language data reflects social and historical inequalities.

Models reproduce stereotypes carried in the data

Problems

1. Modern NLP pipelines rely on pre-trained LLMs
2. Access to sensitive attribute is not automatic



Fairness Through WDM Minimisation⁶

Wasserstein Dependency Measure

Replace KL by a transport distance between joint and product of marginals:

$$I_W = W(p(X;Y); p(X)p(Y)) ;$$

where W is the Wasserstein distance.

 Studying the relation between I_W and fairness metrics.

⁶Leteno T., Perrot M., Laclau C. et al. Fair Classifier via Transferable Representations, JMLR 2025

Fairness Through WDM Minimisation⁶

Wasserstein Dependency Measure

Replace KL by a transport distance between joint and product of marginals:

$$I_W = W(p(X;Y); p(X)p(Y)) ;$$

where W is the Wasserstein distance.

 Studying the relation between I_W and fairness metrics.

We have shown that:

! $I_W(\hat{Y}; A)$ is a linear combination of DP

! $I_W((\hat{Y} = Y) | Y = y; A | Y = y)$ is a linear combination of EO

⁶Leteno T., Perrot M., Laclau C. et al. Fair Classifier via Transferable Representations, JMLR 2025

WDM Regularization with Predicted Sensitive Attribute

Objective

$$\min L(Y; \gamma(X;)) + \lambda_w \hat{Y}; A$$

7 Requires A at train and test time.

WDM Regularization with Predicted Sensitive Attribute

Objective

$$\min L(Y, \gamma(X; \cdot)) + I_W(\hat{Y}, A)$$

7 Requires A at train and test time.

☀ Predict \hat{A} from X to remove dependence on true A .

$$\hat{A} = h_A(\text{Enc}(X)) \quad \text{i.e.} \quad X \xrightarrow{\text{Enc}} Z \xrightarrow{h_A} \hat{A}$$

Lemma 0.2: Approximation of A

We have that:

$$I_W(\hat{Y}, A) = I_W(\hat{Y}, \hat{A}) + 2^p \bar{2}P(A = \hat{A})$$

From Discrete to Continuous Space

Objective

$$\min L(Y, \gamma(X; \cdot)) + I_W(\hat{Y}, \hat{A})$$

7 Argmax operator is non differentiable

From Discrete to Continuous Space

Objective

$$\min L(Y, \gamma(X; \theta)) + I_W(\hat{Y}, \hat{A})$$

7 Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} = Z_y; \hat{A} = Z_a \text{ i.e. } I_W(\hat{Y}, \hat{A}) = I_W(Z_y, Z_a)$$

From Discrete to Continuous Space

Objective

$$\min L(Y, y(X; \theta)) + I_W(\hat{Y}, \hat{A})$$

7 Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} = Z_y; \hat{A} = Z_a \text{ i.e. } I_W(\hat{Y}, \hat{A}) = I_W(Z_y, Z_a)$$

Are the guarantees holding tight? Let's look at the bound!

From Discrete to Continuous Space

Objective

$$\min L(Y, y(X; \cdot)) + I_W(\hat{Y}, \hat{A})$$

7 Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} = Z_y; \hat{A} = Z_a \text{ i.e. } I_W(\hat{Y}, \hat{A}) = I_W(Z_y, Z_a)$$

Are the guarantees holding tight? Let's look at the bound!



From Discrete to Continuous Space

Objective

$$\min L(Y, \gamma(X; \cdot)) + I_W(\hat{Y}, \hat{A})$$

7 Argmax operator is non differentiable

☀ Ensure independence in the latent space

$$\hat{Y} = Z_y; \hat{A} = Z_a \text{ i.e. } I_W(\hat{Y}, \hat{A}) = I_W(Z_y, Z_a)$$

Are the guarantees holding tight? Let's look at the bound!

Forget the equations: it only depends on two intuitive quantities.

- ! **Confidence of predictions:** how clearly one class wins over the others.
- ! **Fraction of uncertain cases:** samples fall below the confidence margin.

Next problem: generative models

Next problem: generative models

Next problem: generative models

Élise, the Nurse

Élise learned the rhythm of the ward the way others learn music.

She could tell, just by the sound of footsteps, who was anxious and who was healing.

One night, she sat a little longer by an elderly patient who couldn't sleep, listening to stories no chart would ever record.

In the morning, she left quietly, knowing care was often invisible but never small.

That was enough for her.

Dr. Martin, the Doctor

Dr. Martin arrived early, reviewing scans while the hospital was still half-asleep.

He liked the certainty of diagnoses, the moment when confusion narrowed into clarity.

But what stayed with him were not the cases he solved, but the conversations he had to slow down for.




Each decision carried weight, and he felt it deeply, even when he hid it well.

By the end of the day, he reminded himself that medicine was as much about people as it was about answers.

Context is everything

What's on your mind today?

+ Who is more at risk of having prostate cancer? John or Jenny?

Outline for today

- ▶ What is fairness?
- ▶ Fair Link Prediction
- ▶ Fair Document Classification
- ▶ **Concluding Remarks**

Fairness is not a plug-in constraint that can be transferred unchanged across learning settings.

Fairness is not a plug-in constraint that can be transferred unchanged across learning settings.

Reasons

A fairness notion is defined relative to:

- ▶ the object being predicted label, ranking, edges, text etc.
- ▶ the intervention point outputs vs representations vs structure.
- ▶ the data-generating system static vs feedback-driven.

Consequence

The question is not “which metric is best?”, but

Which notion matches the system and the harm we aim to prevent?

Any questions ?

