# Multilingual Knowledge Graph Bias Analysis

**Saurav Chennuri**
Department of Computer Science
Boston University
Boston, MA 02215
saurav07@bu.edu

**Zilu Tang**
Department of Computer Science
Boston University
Boston, MA 02215
zilutang@bu.edu

## Abstract

Different cultures exhibit different social biases such as gender inequality. With increasing popularity of multi-lingual knowledge graphs (KG), we are interested in investigating the existence of such differences in biases through analyzing their corresponding KG embeddings. By using the latest graph embedding methods, we attempt to identify the differences in biases in KGs in different languages. Our implementation can be found at [1]

## 1 Introduction

In recent years, large KGs such as DBpedia [2], FreeBase [3], and ConceptNet [4] has attracted increasing interest by the natural language processing (NLP) community. With efforts to improve multilingual knowledge bases, these projects are providing more coverage to culture-specific entities, and improving general concept coverage in low resource languages.

Knowledge bases, however, just like any language models, are prone to inherit social biases from their data sources [Fisher et al., 2019, Blodgett et al., 2020]. If left untreated, such biases will propagated to KG embeddings, which is wildly used in downstream NLP tasks. An interesting question arises as to whether such biases exist uniformly across different KGs in different languages. Since social biases (one example being gender equality) are known to differ between cultures [Group, 2019], one would expect such differences be encoded in the knowledge graphs which represent its "culture".

Recently Vrandecic [2020] proposed an architecture overview for multilingual Wikipedia and how high language resource Wikipedia pages and their low resource counterparts can compensate each other in covering missing knowledge and resolving conflicting topics. Topics such as bias detection and mitigation that we are interested in this paper can be seen as a special case of conflict resolution, where different KGs differ in their biases towards certain topics.

In this project we would like to investigate social biases that exist in different language KGs through hyperbolic embedding methods. Specifically we will implement multi-relational poincaré embedding (MuRP) as well as its euclidean counterpart multi-relational euclidean embedding (MuRE) [Balazevic et al., 2019] and use it as a probing model to investigate bias following the metric defined by Fisher et al. [2019]. More specifically, we would like to investigate the following ideas:

1. How does hyperbolic space change how bias is encoded in the embedding space compared to euclidean space? Does it mitigate or emphasize bias?

2. How are social biases encoded differently in different language knowledge graphs?

---

[1] https://github.com/PootieT/multilingual-knowledge-graph-bias
[2] https://dbpedia.org
[3] https://developers.google.com/freebase
[4] https://conceptnet.io/

## 2 Methods

To understand the advantage of multi-relational embedding, we will need to implement both MuRP and MuRE [Balazevic et al., 2019]. We will also need to implement the bias detection methods in KG introduced by Fisher et al. [2019]. Below sections introduces each ideas in detail.

### 2.1 Multi-relational Poincaré Graph Embedding

Poincaré sphere model of hyperbolic space is one way of representing hyperbolic geometry. Hyperbolic geometry are spaces of constant negative curvature, and they are one form of Riemann manifolds. To understand the poincaré sphere model better, we need to understand about poincaré disk model of hyperbolic space representation. Given this intuition of poincaré disk, we can now move on to understand poincaré sphere which is a 3 dimensional realization of poincaré disk.

In poincaré disk model of hyperbolic space, the embeddings and lines on the hyperbolic space are mapped onto a circular disk in hyperbolic space. On this poincaré disk, the distance between each embedding increases exponentially as we move radially away from the center. This property of poincaré disk is specially useful in providing enough space for embedding heirarchical data that increases exponentially with depth from the root node given a branching factor. So, we use this model of hyperbolic space to embed our data.

There are state of the art models for embedding multi-relational hierarchical data for euclidean spaces and most of them use dot-product as a scoring measure. But there is no clear correspondence to the Euclidean inner product in the hyperbolic space. The Euclidean inner product can be expressed as a function of Euclidean distance and norms as shown in equation 1, i.e.

$$\langle x, y \rangle = \frac{1}{2} \Big( - d_E(x,y)^2 + \|x\|^2 + \|y\|^2 \Big) \tag{1}$$

where x,y are relation adjusted embeddings depending on the geometric space where we are embedding our data.

For poincaré embeddings, the inner product estimation and the relation adjusted embeddings are shown as below

$$\langle x, y \rangle = \frac{1}{2} \Big( - d_E(h_s^{(r)}, h_o^{(r)})^2 + \|x\|^2 + \|y\|^2 \Big) \tag{2}$$

and the poincaré score for relation adjusted embedding for poincaré space as shown as below in equation 3,

$$\langle h_s^r, h_s^r \rangle = \frac{1}{2} \Big( - d_B(exp_0^c(Rlog_0^c(h_s)), h_o \oplus_c r_h)^2 + b_s + b_o \Big) \tag{3}$$

$b_s$ and $b_o$ are bias terms for the tail and head entities. $d_E$ and $d_B$ are distance metrics for euclidean and poincaré spaces as in equations 4 and 5

$$d_E(x,y) = \|x - y\| \tag{4}$$

$$d_B(x,y) = \frac{2}{\sqrt{c}} \tanh^{-1}(\sqrt{c}\| - x \oplus_c y\|) \tag{5}$$

Using this inner product definition in equation 3 as a scoring metric in our poincaré sphere model of hyperbolic space, we embed the entities in our knowledge graph.

### 2.1.1 Visualizing poincaré disk and sphere model

A poincaré disk can be visualized as a circular 2-dimensional disk but in hyperbolic space. The curves and points on the negative curved hyperbolic surface can be mapped to a circular disk and visualized on the disk. It gives an appearance of operating on a disk, but in actuality, we are operating in the hyperbolic geometric space. For visualization, please look at figures 1 and 2
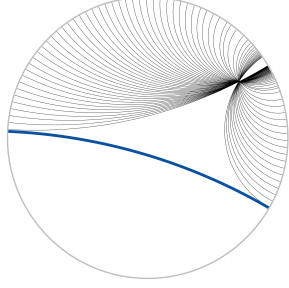
Figure 1: poincaré disk model of hyperbolic space representation. It can be considered as a 2-dimensional interpretation of poincaré sphere
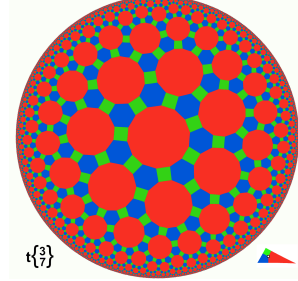


Figure 2: Heirarchical data representation in poincaré disk model

### 2.1.2 Exponential and logarithmic maps

Before we find the poincaré scores, we needed to adjust the hyperbolic entities along with the relations.

We rotate the subject entity in the poincaré sphere and translate the object entity along the direction of relation via möbius addition

To map the lines and embedding of entities onto the poincaré sphere, we map the hyperbolic entity onto the tangent of the poincaré sphere at **0** by logarithmic map, and then align that mapped entity on the tangent along the direction of the relation by multiplying it with a diagonal matrix **M**. After aligning the entity along the relation on the tangent plane, we map it back to the poincaré sphere with exponential map.

This way, we map the hyperbolic subject entities onto the poincaré sphere adjusted to the relation

the exponential and logarithmic maps are defined as follows

$$exp_x^c(v) = x \oplus_c \left( tanh\left(\sqrt{c}\frac{\lambda_x^c\|v\|}{2}\right)\frac{v}{\sqrt{c}\|v\|} \right) \tag{6}$$

$$log_x^c(y) = \frac{2}{\sqrt{c}\lambda_x^c}tanh^{-1}(\sqrt{c}\| - x \oplus_c y\|)\frac{-x \oplus_c y}{\| - x \oplus_c y\|} \tag{7}$$

$$x \oplus_c y = \frac{(1 + 2c\langle \mathbf{x}, \mathbf{y}\rangle + c\|y\|^2)\mathbf{x} + (1 - c\|x\|^2)y)}{1 + 2c\langle \mathbf{x}, \mathbf{y}\rangle + c^2\|x\|^2\|y\|^2} \tag{8}$$

The equations for exponential map, logarithmic map and möbius addition are given in the equations 6, 7 and 8. To better understand visually how these exponential and logarithmic maps work and how we adjust the subject entity according to the relation, please refer to figure 3, figure 4 and figure 5.

### 2.1.3 Riemann Optimization

Since the hyperbolic space we are trying to embed the entities in and train the embeddings, the gradient propagation would also be in the same space. There is an analogy for stochastic gradient descent in hyperbolic space namely "Reimann Stochastic gradient descent"

Since the hyperbolic space are one kind of Riemann manifolds, we use this Riemann optimization method to pass the gradients and train the embeddings. To calculate the gradients in hyperbolic space, we calculate the euclidean gradient $\mathbf{\nabla}_E L$ and multiply it with inverse of Poincaré metric tensor, i.e. $\mathbf{\nabla}_R = 1/(\lambda_\theta^c)^2 \mathbf{\nabla}_E L$. We now use the exponential map $exp_\theta^c$ to project this gradient onto the Poincaré space $\mathbf{\nabla}_R L \in T_\theta B_c^d$ and compute the Riemann update $\theta \leftarrow exp_\theta^c(-\eta\mathbf{\nabla}_R L)$
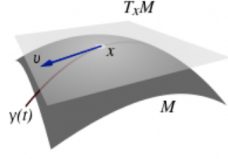
3

Figure 3: Representation of a vector in hyperbolic space
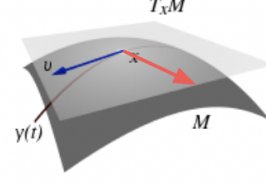


Figure 4: Logarithmic map maps the vector in hyperbolic space to the tangent of the poincaré sphere. Here we map it to a tangent at **0**



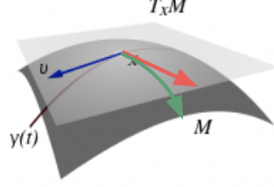Figure 5: exponential map maps the vector on the tangent to the poincaré sphere into the poincaré space

### 2.1.4 Loss function and training the embeddings

The network is implemented in Pytorch. We initialize the input embeddings uniformly at random and pass them through activation functions and operate on them to get the poincaré score output. Ideally if there exists a link between a subject and object entities the score is 1, and if there is no relation, the score is to be zero. Treating this as a binary classification, the Binary Cross Entropy loss function is applied to the outputs of the network which gives the poincaré score

The loss function is given by the equation 9

$$\mathbf{L}(y,p) = -\frac{1}{N}\sum_{i=1}^{N}(y^{(i)}log(p^{(i)}) + (1 - y^{(i)})log(1 - p^{(i)})) \tag{9}$$

### 2.2 Graph Embedding Bias Method

Bias in KG, defined by Fisher et al. [2019] is an unequal probability assigned by a knowledge graph embedding model on a relation, given different values of a sensitive attribute. In our work, we followed previous paper, and focused on investigating the bias in predicting a person's likely profession given different gender of a person (female and male only in our case).

Specifically, the bias is measured with the following approach: for each person entity in the KG, we perturb the embedding to make it more male or female. We then use the model to predict such person's score for having each profession in the KG before and after such perterbation, and use the difference of the score as the bias score for every profession. We repeat such process over all individuals (in large KG we approximate with a subsample of individules), and average the bias score for each profession and report them in a ranked table.

Mathematically, the update of the embedding is defined as:

$$\mathbf{m}(\theta) = \mathbf{g}(e_j, r_{gender}, e_{male}) - \mathbf{g}(e_j, r_{gender}, e_{female}) \tag{10}$$

$$e'_j = e_j + \alpha\frac{\delta\mathbf{m}(\theta)}{\delta e_j} \tag{11}$$

where $e_j$ is the entity embedding of the person $j$. During implementation, we experimented with various way of achieving such result. One way of updating the embedding is providing a pair(s) of triples in a batch $\{e_j, r_{gender}, e_{male}\}$ and $\{e_j, r_{gender}, e_{female}\}$ while assigning the label as 1 and

0 (for biasing towards male, for female, use 0 and 1). Another way, which we found to produce better results, qualitatively, is by providing the previously mentioned triple one by one both with label 1, obtain the gradient with respect to $e_j$, subtract male gradient from female gradient, and update the embedding. During embedding update, for Poincaré, we used Reimannian SGD to the proper scaling factor. After obtaining the score for each person entity, we reload the model weights from checkpoints to ensure the updates to weights do not accumulate across person.

## 3   Datasets

To validate that our graph embedding methods and bias detection methods work, we tested our dataset on FB15K-237[Bordes et al., 2013]. FB15K contains relation triples from Freebase, a large tuple database with structured general human knowledge. It is originally used by Balazevic et al. [2019] as a benchmark for performance comparison with other embedding models.

To train KG embedding on different language KGs, we used different DBpedia KGs that are maintained by a chapter of that specific language. [5]. Specifically we picked English (en), Swedish (sv), and Indonesian (in). Scandanavian cultures are known to have more gender equality. Amongst the languages spoken in Scandanavian countires (Swedish, Norwegian, Finnish, Icelandic), Swedish DBpedia contains the most triples. For similar reasons, we picked Indonesian, out of several other languages spoken in countries with more orthodox social practices (Arabic, Urdo, Hindi)

For all DBpedia datasets, we preproceesd them with the following steps:

1. We first found all variations of the gender entities in a KG and merged different expressions of male and female into either male or female entity. This included converting some surface forms into proper DBpedia entities. For list of all conversions, see data_exploration.py script in the repo. Because there are limited amount of gender triples, we found this necessary in tuning the gender relation to be contextually meaningful.

2. We then extracted all gender related triples and profession related triples (where the relation of the triple is "<http://dbpedia.org/resource/gender" or "<http://dbpedia.org/resource/occupation" (or equivalents in Indonesian and Swedish) and set them aside.

3. We removed any triples with tails that are not a DBpedia resource (mainly text surface form *"The Whale"@en* and literals like *"6"|http://www.w3.org/2001/XMLSchemainteger>*)

4. Because English and Swedish KGs are too big, we filtered them down to a manageable size such that the network can be trained on 10GB GPU. We sampled 10% and 20% of the English and Swedish left-over triples respectively. For Indonesian KG, we did not sub-sample after filtering away literals because the number of entities are at a manageable size.

5. We appended back the gender and profession triples. This way we can keep all gender and profession triples, so all profession and gender entities can more or less keep the original KG's context.

6. We split up all triples into 90/5/5 train, test, dev partitions.

While preprossessing DBpedia data, we found a significant flaw with most of the DBpedia data: most of the people in the KG does not have a gender recorded. Hence in the result table, the sum of the count of female instances of any profession plus the male instance do not add up to the count of the profession in the dataset. We believe this can be addressed in the future by predicting the missing link first and then calculate the bias score. With this fact in mind, it does make one wonder where does the KG learn the concept gender from. We originally abandoned using DBPedia and tried using Wikidata5M [Wang et al., 2019] as an alternative dataset, but sadly found out that it suffers from similar issue, so we went back to DBPedia.

Below is a table documenting the statistics of the filtered datasets:

During model training, we used the following parameter for each datasets we have found:

---

[5]https://docs.google.com/document/d/e/2PACX-1vR7oRoMOkzP5eSLf2vzFPYzJY2BTP-YXdmpmOcczli4GWRZDKq85Ps-DPwbcRJ_xx_UVm4LbargUIay/pub

| Language | unique heads | unique relations | unique tails | Unique Entities |
|---|---|---|---|---|
| FB15K-237 | 13891 | 237 | 13504 | 14541 |
| English | 1540430 | 7486 | 729861 | 2029663 |
| Indonesian | 340865 | 6105 | 412689 | 667714 |
| Swedish | 1578190 | 2624 | 192799 | 1700652 |

Table 1: Entity and Relation Counts

| Language | train triples | dev triples | test triples |
|---|---|---|---|
| FB15K-237 | 272115 | 17535 | 20466 |
| English | 2243938 | 124664 | 124663 |
| Indonesian | 1923392 | 106855 | 106856 |
| Swedish | 2402704 | 133484 | 133484 |

Table 2: Train Validation Test Split Counts

For embedding size we used 40 across all models. For batch size, we fixed it to be as large as we can to speed up the experiment process. In most cases, that is 8192. We begin with a negative sampling size of 50 with learning rate of 50. We found with larger entity size more a larger negative sampling size is necessary to obtain a good evaluation accuracy. For larger datasets (DBpedia KGs), we maximizes this to the largest value that GPU permits, and adjust the learning rate inversely with respect to the batch size. Due to time constraints, all DBpedia models are run for 20 epochs, with two evaluations, one at epoch 10 and epoch 20. During each evaluaiton, we randomly sample a portion of the evaluation set and report the evaluation metrics (such that an evaluation run takes about 5 minutes). For FB15K-237, we ran until the evaluation performance increase become insignificant. Notice, we did not perform extensive hyperparmeter tuning as we are not after the accuracy game. We are more generally interested in the qualitative aspect of our embedding. In general due to more extensive computations performed by Poincaré models, the extra memory requirement during gradient calculation forces us to use a smaller batch size (and hence smaller learning rate as well).

Below is a table of standard link prediction evaluation metrics obtained after training the different models. Hit-k calculates the accuracy of predicting the correct tail within the top-K predictions. Mean rank is the average rank of the correct tail within the final prediction. Mean reciprocal rank shows the average of inverse of the rank of the correct tail in prediction. We included the results from paper as well for a base comparison. We think the reason the author were able to obtain higher accuracy is because they trained for many more epochs with smaller batch size, which we did not have the time for.

As seen from the result, given similar computational resource, within DBpedia KGs, the more relations in a KG, the lower the performance. Swedish link performance results exceeds the other with much less relations and relatively large amount of entities.

# 4   Experiments

We trained MuRP and MuRE on all of the datasets. For each of the dataset, we calculated bias for each of the occupations across all human entities. We display the top male or female biased occupations given the bias score. For KG with larger sets of human entities, we calculate over a subset of them. For English KG, we selected only those professions with more than 10 occurances in the training. We also filtered the people to only those who have more than 5 triples about them. This ensures we have a good quality set of individuals and profession to evaluate over with.

Here are some basic statistics on the dataset:

| Language | Model | Negative Samples | Batch Size | Learning Rate |
|---|---|---|---|---|
| FB15K-237 | Euclidean | 50 | 8192 | 50 |
| FB15K-237 | Poincaré | 50 | 8192 | 50 |
| English | Euclidean | 200 | 4096 | 100 |
| English | Poincaré | 200 | 2048 | 100 |
| Indonesian | Euclidean | 200 | 8192 | 200 |
| Indonesian | Poincaré | 200 | 4096 | 100 |
| Swedish | Euclidean | 200 | 8192 | 200 |
| Swedish | Poincaré | 200 | 4096 | 100 |

Table 3: Model Training Hyperparameters

| Language | Model | Hit-1 | Hit-3 | Hit-10 | Mean Reciprocal Rank | Mean Rank |
|---|---|---|---|---|---|---|
| FB15K-237 (Paper) | Euclidean | .227 | .346 | .493 | .315 | |
| FB15K-237 (Paper) | Poincaré | .235 | .356 | .506 | .324 | |
| FB15K-237 | Euclidean | .176 | .266 | .405 | .251 | 266 |
| FB15K-237 | Poincaré | .172 | .262 | .399 | .247 | 320 |
| English | Euclidean | .050 | .095 | .164 | .087 | 319569 |
| English | Poincaré | .034 | .079 | .143 | .071 | 336469 |
| Indonesian | Euclidean | .197 | .301 | .382 | .262 | 62839 |
| Indonesian | Poincaré | .220 | .308 | .386 | .278 | 64915 |
| Swedish | Euclidean | .339 | .463 | .577 | .421 | 44408 |
| Swedish | Poincaré | .343 | .477 | .612 | .434 | 70272 |

Table 4: Link Prediction Performance

# 5 Results

For readability sake, we only included four tables in the main results section: MuRP and MuRP results on FB15K-237, top male and female biased professions. For results in DBpedia, please refer to the Appendix. Discussions will focus on all of the results (including those in the appendix).

# 6 Discussion

## 6.1 Poincaré bias vs. Euclidean bias

We tried to investigate and compare the biases in both the Euclidean and Poincaré spaces. From our observations on FB15K, English-dbpedia, Indonesian-dbpedia, swedish-dbpedia datasets, we observed that the biases in Euclidean space made more sense than those in poincaré especially in the kinds of language data where the information regarding the number of males and females is not available at large. We knew that poincaré space is better in terms of embedding heirarchical datasets, in comparison to euclidean, but one interesting point to investigate would be how both of these embeddings encode the semantic information that actually exist in real world.

## 6.2 Cross-cultural comparison

Give or take with the quality of our embedding, we can see there are a few unique example professions that stood out in the top male/female biased professions in each of the three KGs. In Indonesia, one of the top male biased profession Caliph (the chief Muslim civil and religious ruler) and Bishop. In addition to religous positions, many included government positions (Member of Pariament, Police, Regent, Parliamentarian, Supreme court, governor, Empress) as well as art and entertainment industry (Actress, Actor, Commedian, Dancer, Songwriter, Screenwriter, Musician, film producer , etc). Top female biased professions in Indonesia include beauty contest, and surprisingly a lot of tech companies (Google, Samsung, Youtube).

As for Swedish results, for top male biased professions, we have mostly jobs in art and entertainment industry (songwriter, singer, musician, rapper, piano, actor) as well as a few leadership positions

| dataset | humans | human subsample fraction | professions | gender triples |
|---|---|---|---|---|
| FB15K-237 | 4532 | 1.0 | 149 | 6093 |
| English (unfiltered) | 67801 | | 7624 | 6093 |
| English | 1075 | 1.0 | 813 | 6093 |
| Indonesian | 13894 | 0.1 | 1760 | 757 |
| Swedish | 2063 | 1.0 | 460 | 1639 |

Table 5: Train Validation Test Split Counts

| Rank | Scores | Profession | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|---|---|---|---|---|---|
| 0 | -4.834e-05 | Association football player | 48 | 36 | 1 |
| 1 | -5.46e-05 | Director of photography | 77 | 65 | 0 |
| 2 | -7.163e-05 | Voice Actor | 347 | 227 | 67 |
| 3 | -7.212e-05 | Cartoonist | 37 | 31 | 0 |
| 4 | -7.491e-05 | Film Directors | 609 | 459 | 42 |
| 5 | -8.377e-05 | screen-writer | 872 | 666 | 71 |
| 6 | -8.457e-05 | Playwright | 107 | 82 | 5 |
| 7 | -8.53e-05 | Comedy performer | 282 | 201 | 45 |
| 8 | -8.532e-05 | Politician | 99 | 68 | 7 |
| 9 | -8.577e-05 | television director | 237 | 183 | 20 |
| 10 | -8.682e-05 | Writter | 441 | 315 | 38 |
| 11 | -8.785e-05 | music theater | 8 | 2 | 0 |
| 12 | -8.96e-05 | engineering (skill) | 54 | 0 | 0 |
| 13 | -9.008e-05 | hollywood producer | 864 | 634 | 80 |
| 14 | -9.154e-05 | Auther | 291 | 187 | 53 |
| 15 | -9.358e-05 | Political Sciences | 125 | 1 | 0 |
| 16 | -9.504e-05 | Seiyû | 28 | 13 | 10 |
| 17 | -9.762e-05 | orchestra conductors | 68 | 51 | 0 |
| 18 | -9.953e-05 | Playback singing | 17 | 10 | 5 |
| 19 | -0.00010134 | Series Producer | 576 | 393 | 79 |
| 20 | -0.00010158 | Television actor | 2271 | 1329 | 545 |

Table 6: Poincaré Embedding 40, FB15K-237, Male Biased top 20 professions

(mayor, lieutenant general, general). Cross-country skiing also appeared as one of the top. For top female biased positions, we saw many law, and government related positions such as (journalist, Politician, Military, lawyer, baron, knight, diplomat, Official), as well as other STEM professions like (doctors, enigneer). Qualitatively, Swedish DBpedia results are the most consistent, and it could be due to the good link prediction results.

Lastly, for English results, the top male biased results include fields such as sports (football association, sport coach, American football official, national football league), to entertainment industry (actor, jockey, rapper, songwriter, TV producer, cinematographer reality TV). Top female results include military positions (Army, US air force), entertainment (Actress, composer, record producer, musician, novelist, voice acting in Japan) to communication/appearance focused professions such as beauty pageant, advertising, publishing, and modeling.

In general, the results are very noisy. One could say the results do indicate that women hold more leadership positions in Swedish society, and that men control more of religious and government positions. However, without clear interpretation of where the source of such bias comes from, it is hard to make a clear conclusion.

### 6.3 Influence from count statistics

In almost all graphs, the is a clear pattern that the top male and female biased professions have higher $C_{total}$ than the others. If we look at the whole table, the occurrence decreases and then increases at the bottom ranks. Assuming our bias calculation implementation is correct (no way of verifying because [Fisher et al., 2019] have not open-sourced their code), this metric seems to be influenced a

| Rank | Scores | Profession | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|------|--------|-----------|--------|--------|----------|
| 0 | -8.111e-05 | Seiyû | 28 | 13 | 10 |
| 1 | -8.926e-05 | orchestra conductors | 68 | 51 | 0 |
| 2 | -9.32e-05 | timothy thompson (composer) | 272 | 211 | 11 |
| 3 | -9.408e-05 | Playback singing | 17 | 10 | 5 |
| 4 | -0.00010226 | Film Directors | 609 | 459 | 42 |
| 5 | -0.00010465 | Painoist | 53 | 37 | 4 |
| 6 | -0.00010539 | Songwriting | 380 | 242 | 65 |
| 7 | -0.0001059 | guitarrist | 157 | 115 | 15 |
| 8 | -0.00010746 | lyricists | 82 | 63 | 7 |
| 9 | -0.00010754 | music production (music industry) | 291 | 206 | 35 |
| 10 | -0.00010788 | keyboarder | 65 | 48 | 5 |
| 11 | -0.00011099 | Director of photography | 77 | 65 | 0 |
| 12 | -0.00011149 | Singer songwriter | 286 | 163 | 70 |
| 13 | -0.00011264 | Contrabassist | 30 | 23 | 2 |
| 14 | -0.00011291 | music career | 585 | 418 | 61 |
| 15 | -0.00011354 | Concert organist | 23 | 16 | 0 |
| 16 | -0.00011454 | Sister newspaper | 4 | 1 | 0 |
| 17 | -0.00011514 | Drummers | 30 | 23 | 0 |
| 18 | -0.00011521 | Multiinstrumentalist | 28 | 19 | 4 |
| 19 | -0.00011549 | Political Sciences | 125 | 1 | 0 |
| 20 | -0.00011555 | playwright | 107 | 82 | 5 |

Table 7: Euclidean Embedding 40, FB15K-237, Male Biased top 20 professions

| Rank | Scores | Profession | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|------|--------|-----------|--------|--------|----------|
| 0 | 0.00049628 | makeup artist | 807 | 2 | 0 |
| 1 | 0.00049535 | special effects coordinator | 668 | 2 | 0 |
| 2 | 0.00049266 | Sound editor (filmmaking) | 321 | 1 | 0 |
| 3 | 0.00045964 | storyboard artist | 106 | 2 | 0 |
| 4 | 0.0004478 | communication designer | 101 | 2 | 0 |
| 5 | 0.00044708 | hairdresser | 68 | 2 | 0 |
| 6 | 0.00043116 | Digital director | 175 | 4 | 0 |
| 7 | 0.00039009 | animation direction | 34 | 2 | 0 |
| 8 | 0.00034352 | Black and white artist | 26 | 6 | 0 |
| 9 | 0.00033142 | Not Found | 96 | 16 | 1 |
| 10 | 0.00031555 | animated film director | 114 | 26 | 0 |
| 11 | 0.00028239 | choreographic technique | 6 | 1 | 0 |
| 12 | 0.00024113 | Impresario | 1 | 1 | 0 |
| 13 | 0.0002403 | Head basketball coach | 1 | 1 | 0 |
| 14 | 0.00023955 | geologists | 1 | 1 | 0 |
| 15 | 0.00023926 | socialite | 1 | 0 | 1 |
| 16 | 0.00023715 | Biologist | 1 | 1 | 0 |
| 17 | 0.00023708 | backup vocal | 19 | 3 | 0 |
| 18 | 0.00023619 | event promotion | 1 | 0 | 0 |
| 19 | 0.00023386 | oncological biochemist | 1 | 1 | 0 |
| 20 | 0.00023314 | Motivational speaking | 1 | 1 | 0 |

Table 8: Poincaré Embedding 40, FB15K-237, Female Biased top 20 professions

lot by simple count statistics. Perhaps with more count, a profession entity is more likely connected with gender somehow, hence be influenced by the gender triple more during perturbation. On the other hand, the pattern is not full-proof. There are instances of professions with high count but is ranked lower and professions with low counts ranked higher. We hypothesize that these are the professions which contain "proportionally" more bias. For example, in **??**, "animated film director" has 114 counts while four of the higher ranked professions ("hairdresser", "animation direction",

| Rank | Scores | Profession | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|---|---|---|---|---|---|
| 0 | 0.00016167 | makeup artist | 807 | 2 | 0 |
| 1 | 0.00016148 | special effects coordinator | 668 | 2 | 0 |
| 2 | 0.00016107 | Sound editor (filmmaking) | 321 | 1 | 0 |
| 3 | 0.00015546 | Hairdresser | 68 | 2 | 0 |
| 4 | 0.00015535 | storyboard artist | 106 | 2 | 0 |
| 5 | 0.00015376 | communication designer | 101 | 2 | 0 |
| 6 | 0.00015147 | Digital director | 175 | 4 | 0 |
| 7 | 0.00014739 | animation direction | 34 | 2 | 0 |
| 8 | 0.00014349 | Not Found | 96 | 16 | 1 |
| 9 | 0.00014293 | Black and white artist | 26 | 6 | 0 |
| 10 | 0.00014178 | animated film director | 114 | 26 | 0 |
| 11 | 0.00013777 | female model | 150 | 32 | 90 |
| 12 | 0.00013702 | stunt person | 11 | 6 | 0 |
| 13 | 0.00013384 | choreographic technique | 6 | 1 | 0 |
| 14 | 0.0001313 | socialite | 1 | 0 | 1 |
| 15 | 0.00013099 | designer | 20 | 5 | 5 |
| 16 | 0.00013056 | event promotion | 1 | 0 | 0 |
| 17 | 0.0001305 | Head basketball coach | 1 | 1 | 0 |
| 18 | 0.00013036 | super-models | 4 | 0 | 4 |
| 19 | 0.00013035 | geologists | 1 | 1 | 0 |
| 20 | 0.00013005 | Biologist | 1 | 1 | 0 |

Table 9: Euclidean Embedding 40, FB15K-237, Female Biased top 20 professions

"Black and white artist", "Not Found") have lower count than its count. With a simple google search, it's not hard to infer that "animated film director" is a male-dominated field.



Table 10: Result of searching "animated film director" in Google

## 6.4 Challenges with missing gender triples

In all of the graphs, we noticed significant amount of gender triples missing. Proportionally, FB15K-237 had the most gender triples and we were actually able to see them in the tables. However, there appears to be a lot of erroneous links as well. For instance there are 32 "female models" that are of gender "male" in Table 9. Given missing information, it is hard to tell what are the sources of triples which influenced the entity embeddings. An interesting future direction is to use techniques like influence function Han et al. [2020] to provide more interpretability to node embeddings. On the other hand, it would be worth to investigate what properties in the knowledge graph led to such embedding. If a profession is biased towards male, we can calculate the min-cut, shortest path, number of paths of length K, etc between male and the profession entity on the graph where we treat all relations as edges. Given different relations, we can also try to calculate paths between the two nodes by allowing only subset of relations as possible edges.

## 6.5 Challenges with Subsampling KG

One of the main challenge when training a KG embedding is the size of the KG. Given limited hardware 10G GPU memory, the maximum entity size is around 2.4 M nodes with node embedding dimension 40. Due to such constraint, we needed to subsample the larger KG graphs. During the

process of sampling, we also want to keep the "structure" of the graph relatively intact. We want to maximize the number of triples per relation or entity, to avoid infrequent entities or relations not yielding good embedding due to lack of training data. We also want to maintain somewhat of a diverse set of relations so embeddings of entities benefit from various types of relations as well as connected neighbors. There are a few ways we thought of sampling but all have their downsides:

- Pick top few most frequent relations or entities. This way we maximize the "quality" of the top few entities / relations' embedding. Unfortunately, we found most of these entities / relations follow the power distribution with long tail. If we want 10% of the original data, we might only get one entity left in the training data.

- Perform weighted sampling given the count of entity / relation occurrence or sum of tail and head occurrence in dataset. This is more or less a better balance between two objectives, and we can tune the preference by using $log$ of count or exponentiating the count to a certain power.

- Randomly sample triples in KG. This is the method we ended up with. However, more questions came up while we were thinking in this direction: what is the best trade-off when maintaining the "quality" of a KG by subsampling relations and entities? How important is entity diversity? How important do the few most frequent entities and relations play in defining the "landscape" of the knowledge graph?

## 6.6  Future Directions

The following mentioned ideas could be explored by building on our work for future

- Check the bias for other sensitive attributes like "Race", "community", "Ethnicity", "Religion" ,with other factors like "occupation" etc.

- Biases difference between Poincaré and euclidean model (in hierarchical relationship specifically). One interesting future direction is to investigate in KGs with more hierarchies, and look for relations that are hierarchical in nature. If Poincare represents hierarchical graph better, does it mean it encode the "bias" better as well?

- Identify source of bias given graph properties. We could possible apply the following methods like Shortest path,  paths, Min-cut between sensitive attributes and other properties in pursuit of finding the source of bias

- How to sparsify KG graph so it maintain the hierarchical structure well enough for link prediction

- Bias mitigation on top of knowledge graph embeddings

## References

I. Balazevic, C. Allen, and T. Hospedales. Multi-relational poincaré graph embeddings. *Advances in Neural Information Processing Systems*, 32:4463–4473, 2019.

S. L. Blodgett, S. Barocas, H. Daumé III, and H. Wallach. Language (technology) is power: A critical survey of" bias" in nlp. *arXiv preprint arXiv:2005.14050*, 2020.

A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko. Translating embeddings for modeling multi-relational data. In C. J. C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. URL https://proceedings.neurips.cc/paper/2013/file/1cecc7a77928ca8133fa24680a88d2f9-Paper.pdf.

J. Fisher, D. Palfrey, C. Christodoulopoulos, and A. Mittal. Measuring social bias in knowledge graph embeddings. *arXiv preprint arXiv:1912.02761*, 2019.

W. B. Group. *Women, business and the law 2019: A decade of reform*. World Bank, 2019.

X. Han, B. C. Wallace, and Y. Tsvetkov. Explaining black box predictions and unveiling data artifacts through influence functions. *arXiv preprint arXiv:2005.06676*, 2020.

D. Vrandecic. Architecture for a multilingual wikipedia. 2020.

X. Wang, T. Gao, Z. Zhu, Z. Liu, J. Li, and J. Tang. Kepler: A unified model for knowledge embedding and pre-trained language representation. *arXiv preprint arXiv:1911.06136*, 2019.

# A  Appendix

## A.1  Indonesia DBpedia Results

| Rank | Profession | Scores | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|---|---|---|---|---|---|
| 0 | Empress | 0.0020534 | 84 | 0 | 0 |
| 1 | Caliph | 0.00191078 | 75 | 0 | 0 |
| 2 | Mother_Country_United States_United States | 0.00168522 | 88 | 0 | 0 |
| 3 | Bishop | 0.00167088 | 91 | 0 | 0 |
| 4 | Novel | 0.00140156 | 147 | 0 | 0 |
| 5 | Member_Parliament_Europe | 0.00115505 | 42 | 0 | 0 |
| 6 | Rock_and_roll | 0.00114182 | 138 | 0 | 0 |
| 7 | Wife_Vice_President_United States_United States | 0.00100457 | 38 | 0 | 0 |
| 8 | Police_State_Republic_Indonesia | 0.00097306 | 251 | 0 | 0 |
| 9 | Classic_music | 0.00096806 | 94 | 0 | 0 |
| 10 | Regent | 0.00096412 | 418 | 0 | 0 |
| 11 | Dangdut | 0.00095131 | 295 | 0 | 0 |
| 12 | Rap | 0.00083163 | 95 | 0 | 0 |
| 13 | Rock | 0.00078534 | 551 | 0 | 0 |
| 14 | Jazz | 0.00071153 | 586 | 0 | 0 |
| 15 | Parliamentarian | 0.0006188 | 125 | 0 | 0 |
| 16 | entertainment | 0.00053724 | 107 | 0 | 0 |
| 17 | Internet | 0.00045979 | 73 | 0 | 0 |
| 18 | Pharmacy | 0.00036151 | 50 | 0 | 0 |
| 19 | Supreme Court_Republic_Indonesia | 0.00018654 | 69 | 0 | 0 |
| 20 | governor | 0.00011319 | 92 | 0 | 0 |

Table 11: Poincaré Embedding 40, Indonesian DBpedia, Male Biased top 20 professions

## A.2  Swedish DBpedia Results

## A.3  English DBpedia Results

| Rank | Profession | Scores | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|---|---|---|---|---|---|
| 0 | Trans_Media | 0.00479248 | 67 | 0 | 0 |
| 1 | Cosmonaut | 0.00449267 | 63 | 0 | 0 |
| 2 | Book | 0.00428434 | 66 | 0 | 0 |
| 3 | Song | 0.00424034 | 54 | 0 | 0 |
| 4 | Fighter | 0.00412852 | 32 | 0 | 0 |
| 5 | Organization_nonprofit | 0.00412557 | 38 | 0 | 0 |
| 6 | Beauty contest | 0.00412078 | 105 | 0 | 0 |
| 7 | vice-regent | 0.00382302 | 205 | 0 | 0 |
| 8 | Army | 0.00358525 | 31 | 0 | 0 |
| 9 | Era_Commander_War | 0.00356383 | 12 | 0 | 0 |
| 10 | Google | 0.00355191 | 108 | 0 | 0 |
| 11 | Private | 0.0034907 | 108 | 0 | 0 |
| 12 | Marines | 0.00338265 | 13 | 0 | 0 |
| 13 | Doctor | 0.00337311 | 178 | 0 | 0 |
| 14 | Schutzstaffel | 0.00335675 | 9 | 0 | 0 |
| 15 | Bachelor of Law | 0.00326354 | 79 | 0 | 0 |
| 16 | NGO | 0.00322126 | 7 | 0 | 0 |
| 17 | Samsung | 0.00321324 | 26 | 0 | 0 |
| 18 | MD_Entertainment | 0.00321127 | 305 | 0 | 0 |
| 19 | Referee_(soccer_soccer) | 0.00320898 | 66 | 0 | 0 |
| 20 | TNI | 0.00320402 | 193 | 0 | 0 |

Table 12: Poincaré Embedding 40, Indonesian DBpedia, Female Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|---|---|---|---|---|---|
| 0 | Regent | 6.026e-05 | 418 | 0 | 0 |
| 1 | Mother_Country_United States_United States | 3.57e-05 | 88 | 0 | 0 |
| 2 | Empress | 9.79e-06 | 84 | 0 | 0 |
| 3 | Model | -3.345e-05 | 713 | 0 | 0 |
| 4 | Actress | -4.842e-05 | 2281 | 0 | 0 |
| 5 | Model_(Job) | -5.487e-05 | 401 | 0 | 0 |
| 6 | Master of Ceremony | -5.858e-05 | 363 | 0 | 0 |
| 7 | Producer | -6.809e-05 | 256 | 0 | 0 |
| 8 | Cast | -7.049e-05 | 654 | 4 | 1 |
| 9 | Presenter | -7.734e-05 | 284 | 0 | 0 |
| 10 | Comedian | -8.469e-05 | 194 | 0 | 0 |
| 11 | Caliph | -8.683e-05 | 75 | 0 | 0 |
| 12 | Actor | -8.89e-05 | 2364 | 0 | 5 |
| 13 | Dancer | -9.023e-05 | 294 | 1 | 0 |
| 14 | Songwriter | -9.123e-05 | 374 | 0 | 0 |
| 15 | Screenwriter | -9.646e-05 | 222 | 0 | 0 |
| 16 | Singer | -9.681e-05 | 2675 | 5 | 2 |
| 17 | musician | -0.00010193 | 359 | 0 | 0 |
| 18 | Producer_film | -0.00010885 | 157 | 0 | 1 |
| 19 | Vocals | -0.00011022 | 870 | 1 | 1 |
| 20 | Businessman | -0.00011273 | 459 | 0 | 0 |

Table 13: Euclidean Embedding 40, Indonesian DBpedia, Male Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ | $C_{male}$ | $C_{female}$ |
|---|---|---|---|---|---|
| 0 | Cosmonaut | 0.00132466 | 63 | 0 | 0 |
| 1 | Beauty contest | 0.00128166 | 105 | 0 | 0 |
| 2 | Song | 0.00122522 | 54 | 0 | 0 |
| 3 | Trans_Media | 0.0012137 | 67 | 0 | 0 |
| 4 | Referee_(soccer_soccer) | 0.00113927 | 66 | 0 | 0 |
| 5 | Fighter | 0.00112451 | 32 | 0 | 0 |
| 6 | Organization_nonprofit | 0.00110462 | 38 | 0 | 0 |
| 7 | YouTube | 0.00105859 | 67 | 0 | 0 |
| 8 | Private | 0.00105121 | 108 | 0 | 0 |
| 9 | Antv | 0.00104173 | 82 | 0 | 0 |
| 10 | EMI | 0.00103991 | 154 | 0 | 0 |
| 11 | MD_Entertainment | 0.00101431 | 305 | 0 | 0 |
| 12 | Army | 0.00101251 | 31 | 0 | 0 |
| 13 | South Kalimantan | 0.0009661 | 194 | 0 | 0 |
| 14 | Google | 0.0008913 | 108 | 0 | 0 |
| 15 | Book | 0.00084801 | 66 | 0 | 0 |
| 16 | Television | 0.00084777 | 105 | 1 | 0 |
| 17 | Catholic_Roman | 0.00084057 | 1997 | 3 | 1 |
| 18 | Samsung | 0.0008253 | 26 | 0 | 0 |
| 19 | Indonesia | 0.00080967 | 13117 | 3 | 0 |
| 20 | Islam | 0.0008077 | 6198 | 7 | 4 |

Table 14: Euclidean Embedding 40, Indonesian DBpedia, Female Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ |
|---|---|---|---|
| 0 | Songwriter | 0.00183121 | 482 |
| 1 | Singer | 0.00162604 | 525 |
| 2 | Musician | 0.00162394 | 429 |
| 3 | Mayor | 0.00156511 | 293 |
| 4 | Astronomy | 0.00075089 | 70 |
| 5 | Composer | 0.00059657 | 70 |
| 6 | Actor | 0.00058228 | 266 |
| 7 | Singer | 0.00056863 | 45 |
| 8 | Singer-songwriter | 0.00054025 | 44 |
| 9 | Piano | 0.00049689 | 213 |
| 10 | Painting art | 0.00041926 | 40 |
| 11 | Dancer | 0.0003884 | 41 |
| 12 | Rapper | 0.00037003 | 19 |
| 13 | Law degree | 0.00034568 | 11 |
| 14 | Conductor | 0.00033878 | 21 |
| 15 | Los_Angeles_Kings | 0.00033103 | 60 |
| 16 | songwriter | 0.00032011 | 13 |
| 17 | Youth literature | 0.00031553 | 27 |
| 18 | Sculptor | 0.00031297 | 21 |
| 19 | Cross-country skiing | 0.00030766 | 52 |
| 20 | Lieutenant general | 0.00030088 | 42 |

Table 15: Euclidean Embedding 40, Swedish DBpedia, Male Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ | $C_{female}$ |
|---|---|---|---|---|
| 0 | Stockholm | 0.00017633 | 1187 | 0 |
| 1 | Journalist | 9.765e-05 | 530 | 1 |
| 2 | Politician | 6.57e-05 | 363 | 0 |
| 3 | Author | 5.316e-05 | 1827 | 0 |
| 4 | Military | 4.739e-05 | 280 | 0 |
| 5 | Lawyer | 4.649e-05 | 285 | 1 |
| 6 | Lawyer | 2.244e-05 | 267 | 0 |
| 7 | Translator | 1.397e-05 | 216 | 0 |
| 8 | Swedish Church | -1.115e-05 | 640 | 0 |
| 9 | Travkusk | -2.897e-05 | 159 | 0 |
| 10 | Trotting | -3.367e-05 | 143 | 0 |
| 11 | Poet | -5.974e-05 | 184 | 0 |
| 12 | Businessman | -7.509e-05 | 119 | 0 |
| 13 | Diplomat | -8.103e-05 | 106 | 0 |
| 14 | Scriptwriter | -9.757e-05 | 94 | 0 |
| 15 | Teacher | -0.00010664 | 144 | 0 |
| 16 | Doctor | -0.00012358 | 75 | 0 |
| 17 | Playwright | -0.00012988 | 72 | 0 |
| 18 | Series creator | -0.00013092 | 67 | 0 |
| 19 | Engineer | -0.00013653 | 65 | 0 |
| 20 | Official | -0.00014054 | 66 | 0 |

Table 16: Euclidean Embedding 40, Swedish DBpedia, Female Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ |
|---|---|---|---|
| 0 | Songwriter | 0.00482701 | 482 |
| 1 | Mayor | 0.00445378 | 293 |
| 2 | Musician | 0.00436379 | 429 |
| 3 | Singer | 0.00431085 | 525 |
| 4 | Astronomy | 0.00218133 | 70 |
| 5 | Composer | 0.0018701 | 70 |
| 6 | Singer | 0.00179402 | 45 |
| 7 | Piano | 0.00174245 | 213 |
| 8 | Singer-songwriter | 0.00162994 | 44 |
| 9 | Actor | 0.00143001 | 266 |
| 10 | Painting art | 0.00137375 | 40 |
| 11 | Dancer | 0.00117425 | 41 |
| 12 | Law degree | 0.0011027 | 11 |
| 13 | Rapper | 0.00108797 | 19 |
| 14 | Master of Science in Engineering | 0.00098262 | 31 |
| 15 | Nintendo | 0.00098074 | 80 |
| 16 | Conductor | 0.00097519 | 21 |
| 17 | Stockholm | 0.00092926 | 1187 |
| 18 | Sculptor | 0.00090629 | 21 |
| 19 | Lieutenant general | 0.00090212 | 42 |
| 20 | General | 0.00088431 | 36 |

Table 17: Poincaré Embedding 40, Swedish DBpedia, Male Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ | $C_{female}$ |
|------|------------|--------|-------------|--------------|
| 0 | Swedish Church | 0.00010978 | 640 | 0 |
| 1 | Journalist | 2.423e-05 | 530 | 1 |
| 2 | Author | 2.233e-05 | 1827 | 0 |
| 3 | Politician | -6.072e-05 | 363 | 0 |
| 4 | Military | -9.266e-05 | 280 | 0 |
| 5 | Lawyer | -9.622e-05 | 285 | 1 |
| 6 | Trotting | -0.00012097 | 143 | 0 |
| 7 | Lawyer | -0.00014641 | 267 | 0 |
| 8 | Translator | -0.00015073 | 216 | 0 |
| 9 | Baron | -0.0002045 | 25 | 0 |
| 10 | Knight | -0.00020719 | 23 | 0 |
| 11 | Travkusk | -0.0002341 | 159 | 0 |
| 12 | Poet | -0.00029711 | 184 | 0 |
| 13 | Diplomat | -0.00032036 | 106 | 0 |
| 14 | Businessman | -0.00032379 | 119 | 0 |
| 15 | Scriptwriter | -0.00036458 | 94 | 0 |
| 16 | Doctor | -0.00040115 | 75 | 0 |
| 17 | Teacher | -0.00040877 | 144 | 0 |
| 18 | Series creator | -0.0004101 | 67 | 0 |
| 19 | Playwright | -0.00042537 | 72 | 0 |
| 20 | Engineer | -0.00043338 | 65 | 0 |

Table 18: Poincaré Embedding 40, Swedish DBpedia, Female Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ |
|------|------------|--------|-------------|
| 0 | Reality_television | 0.00069947 | 175 |
| 1 | Classical_music | 0.0006744 | 300 |
| 2 | Beauty_pageant | 0.00065037 | 66 |
| 3 | Film_score | 0.00062718 | 96 |
| 4 | Bishop | 0.00057922 | 106 |
| 5 | Talk_show | 0.00053672 | 65 |
| 6 | General | 0.00053497 | 104 |
| 7 | Poetry | 0.00050894 | 116 |
| 8 | Hindustani_classical_music | 0.00050857 | 49 |
| 9 | Documentary_film | 0.00050543 | 55 |
| 10 | Short_story | 0.00050342 | 73 |
| 11 | Nonprofit_organization | 0.00047156 | 146 |
| 12 | Opera | 0.00045918 | 51 |
| 13 | Sheriff | 0.00043167 | 23 |
| 14 | Chamber_music | 0.00041481 | 28 |
| 15 | Literary_criticism | 0.00041106 | 35 |
| 16 | Association_football | 0.00040854 | 113 |
| 17 | Pirate | 0.00040412 | 25 |
| 18 | Doctor_(title) | 0.00040279 | 30 |
| 19 | Public_house | 0.00039895 | 14 |
| 20 | Coach_(sport) | 0.00039812 | 28 |

Table 19: Euclidean Embedding 40, English DBpedia, Male Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ |
|---|---|---|---|
| 0 | Singing | 1.805e-05 | 1310 |
| 1 | Film_director | -6.13e-06 | 3274 |
| 2 | Actress | -7.73e-06 | 3043 |
| 3 | Actor | -8.67e-06 | 5862 |
| 4 | Songwriter | -1.122e-05 | 1765 |
| 5 | Composer | -1.175e-05 | 1712 |
| 6 | Screenwriter | -1.177e-05 | 2480 |
| 7 | Record_producer | -1.206e-05 | 1788 |
| 8 | Musician | -1.766e-05 | 1824 |
| 9 | Author | -1.813e-05 | 1806 |
| 10 | Novelist | -1.914e-05 | 1198 |
| 11 | Film_producer | -2.19e-05 | 1580 |
| 12 | Politician | -2.217e-05 | 4437 |
| 13 | Journalist | -2.272e-05 | 2005 |
| 14 | Lawyer | -2.555e-05 | 3563 |
| 15 | Model_(person) | -2.826e-05 | 872 |
| 16 | Poet | -2.927e-05 | 991 |
| 17 | Teacher | -3.758e-05 | 803 |
| 18 | Television_presenter | -3.847e-05 | 754 |
| 19 | Cinematographer | -3.997e-05 | 669 |
| 20 | Voice_acting_in_Japan | -4.186e-05 | 625 |

Table 20: Euclidean Embedding 40, English DBpedia, Female Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ |
|---|---|---|---|
| 0 | American_football_official | 0.00010667 | 40 |
| 1 | National_Football_League | 9.837e-05 | 60 |
| 2 | Cinematographer | -6.878e-05 | 669 |
| 3 | Film_editor | -0.00012494 | 243 |
| 4 | General | -0.00012996 | 104 |
| 5 | Film_actor | -0.00013622 | 452 |
| 6 | Film_director | -0.00014779 | 3274 |
| 7 | Television_director | -0.00014978 | 268 |
| 8 | Actress | -0.00014995 | 3043 |
| 9 | Composer | -0.00014998 | 1712 |
| 10 | Actor | -0.00015153 | 5862 |
| 11 | Model_(person) | -0.0001528 | 872 |
| 12 | Novelist | -0.00016012 | 1198 |
| 13 | Conducting | -0.00016111 | 299 |
| 14 | Jockey | -0.00016168 | 510 |
| 15 | Rapper | -0.0001617 | 615 |
| 16 | Game_designer | -0.00016205 | 214 |
| 17 | Voice_acting_in_Japan | -0.00016307 | 625 |
| 18 | Lyricist | -0.00016672 | 235 |
| 19 | Songwriter | -0.00016786 | 1765 |
| 20 | Television_producer | -0.00016873 | 531 |

Table 21: Poincaré Embedding 40, English DBpedia, Male Biased top 20 professions

| Rank | Profession | Scores | $C_{total}$ |
|------|------------|--------|-------------|
| 0 | Member_of_Parliament | 0.00086142 | 337 |
| 1 | Bishop | 0.00083277 | 106 |
| 2 | Reality_television | 0.00077934 | 175 |
| 3 | Film_score | 0.00074491 | 96 |
| 4 | Talk_show | 0.00072417 | 65 |
| 5 | Classical_music | 0.00071962 | 300 |
| 6 | Documentary_film | 0.00070873 | 55 |
| 7 | United_States_Air_Force | 0.00070603 | 218 |
| 8 | Hindustani_classical_music | 0.00070291 | 49 |
| 9 | Short_story | 0.00068853 | 73 |
| 10 | Beauty_pageant | 0.00064871 | 66 |
| 11 | Private_equity | 0.00063567 | 32 |
| 12 | Poetry | 0.00062741 | 116 |
| 13 | Opera | 0.00061221 | 51 |
| 14 | Member_of_the_Legislative_Assembly_(India) | 0.00059806 | 271 |
| 15 | Stand-up_comedy | 0.00058327 | 55 |
| 16 | Real_estate | 0.00058226 | 61 |
| 17 | Coach_(sport) | 0.00058005 | 28 |
| 18 | Publishing | 0.00057161 | 52 |
| 19 | Advertising | 0.00054637 | 29 |
| 20 | Army | 0.00053263 | 69 |

Table 22: Poincaré Embedding 40, English DBpedia, Female Biased top 20 professions