vector Representations (Aka Embeddings)

- Approach 1   Dictionary
X

$w_1$ ☐

$w_2$ ☐

...

$w_N$ ☐

output

X

X
Sentence Structure
Semantic Sim. X
Sentiment X
...

You Cannot Piggyback on word Sim. for training the model

- approach 2:

allocate a vector of random numbers to each word

| 1 | 2 | .... | | d |
|---|---|------|---|---|
| 5 | 11 | | 6 | 9 |

X does not Consider Semantic Similarities.

---

Vector Representations (Embeddings)
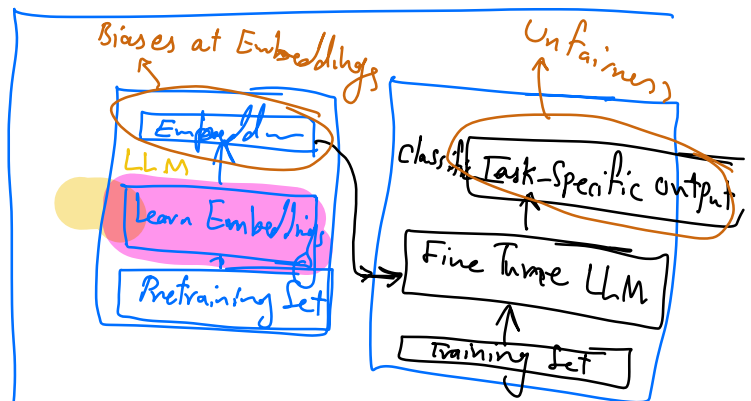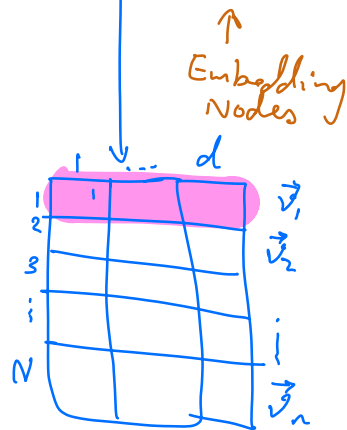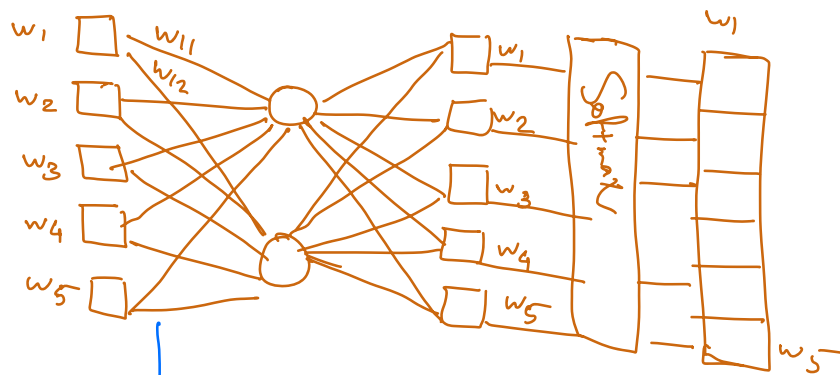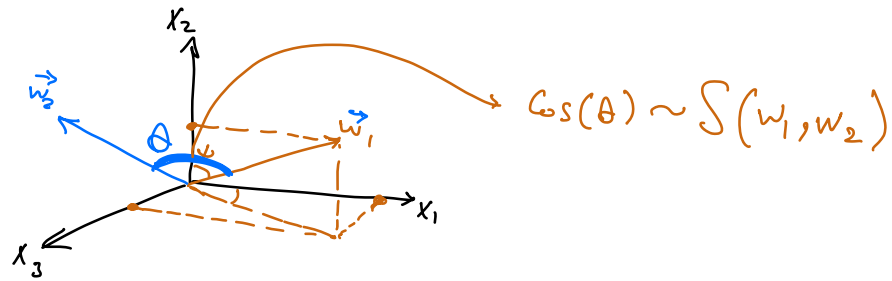
Given a Universe $U = \{u_1, \ldots u_N\}$,

a function $S$, $S(u_i, u_i)$ : Semantic Sim. b/ $u_i, u_j$

$\vec{v_i} = Vec(u_i)$   is a   d-dimensional vector of numbers

s.t.

$$Cos \angle (\vec{v_i}, \vec{v_j}) \sim S(u_i, u_j)$$

$$\cos(\theta) \sim S(w_1, w_2)$$

$w_1$ $w_{11}$ $w_1$ $w_1$
$w_2$ $w_{12}$ $w_2$ Softmax
$w_3$ $w_3$
$w_4$ $w_4$
$w_5$ $w_5$ $w_5$

Embedding Nodes

$1 \cdots d$
$1$
$2$ $\vec{v_1}$
$3$ $\vec{v_2}$
$i$ $i$
$N$ $\vec{v_n}$

Biases at Embeddings

Unfairness

Embedding

LLM

Learn Embeddings

Pretraining Set

Classif. Task-Specific output

Fine Tune LLM

Training Set

Embedding of

Embeddings are biased if the words that should be
nuetral are more similar to one of
the groups.

word

e.g.,

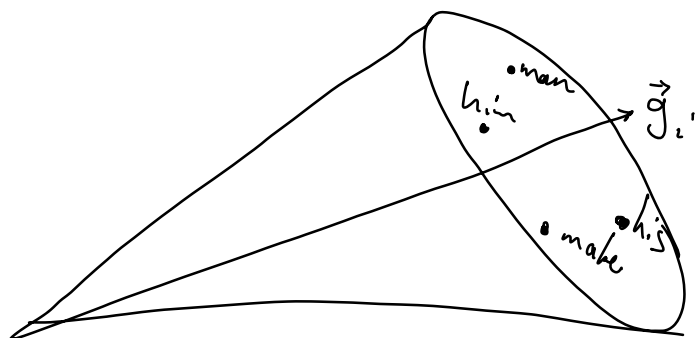She nurse doctor
$\theta_2$
he
$\theta_1$ $\theta_1 < \theta_2$

— Every word semantically representing $g_i$ (male, man, him, his...) should have a small angle $\vec{g_i}$, $\Rightarrow$ Falling inside its cone.

**Assumption** → Every Representative word, is an iid sample from the surface of the cone



$$E\left[\vec{g_i}\right] = \frac{1}{K} \sum \vec{\nu_i}$$

Given a Set of words $\{A, B\}$

$\{$Doctor, Eng ...$\}$ $\rightarrow$

$\{$Nurse, ...$\}$ $\rightarrow$

$$WEAT = \frac{1}{\ell}\left[ \sum_{a_i \in A} Sim(\vec{a_i}, \vec{g_1}) - Sim(\vec{a_j}, \vec{g_2}) \right.$$
$$\left. + \sum_{b_i \in B} Sim(\vec{b_i}, \vec{g_2}) - Sim(\vec{b_i}, g_1) \right]$$

$$Sim(\vec{a}, \vec{g}) = Cos \angle (\vec{a}, \vec{g})$$

$\Downarrow$ Extension To Sentences

Embedding $\mathcal{E}$: Sentece $\rightsquigarrow \vec{v}$

$\quad \hookrightarrow$ Instructor.

measuring bias of a Sentence Embedding

Step 1: Finding Embeddings for groups

$g_1$ : male $\qquad\qquad$ $g_2$ : Female

- Use a Predefined Set of Sentences:

    A) have a high association with $g_i$
    B) " a minimal (to Zero) Extra
        Information.
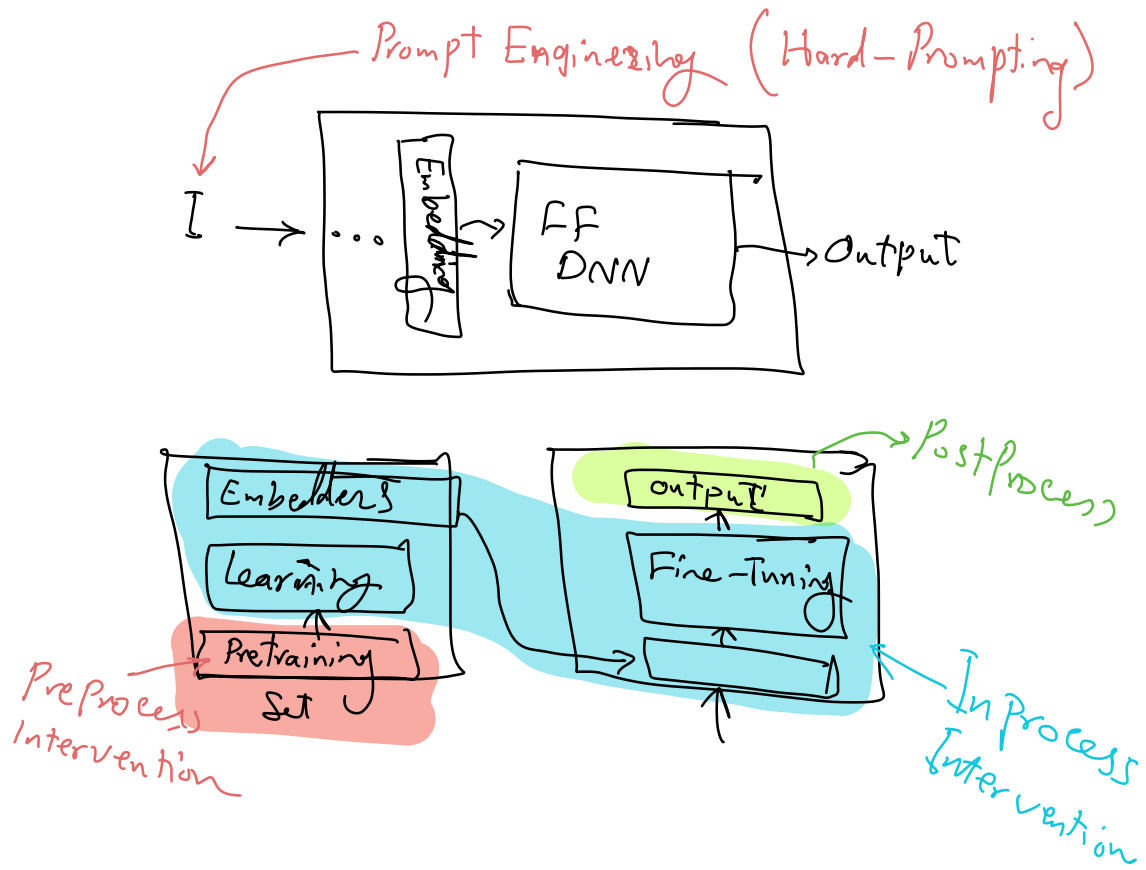
{ "He is a boy", "His gender is Male", ... }

- Assume each Sentence is an IID Sample
  in the "Cone" of $g_i$

- Take the average as the Expected
  Value of $\vec{g}_i$

$$\vec{g}_i = \frac{1}{K} \sum_{i=1}^{K} \vec{v}_i$$

$$SEAT = \frac{1}{\ell} \left( \sum_{a \in A} Sim(a, \vec{g}_1) - Sim(a, \vec{g}_2) \right.$$
$$\left. + \sum_{b \in B} Sim(b, \vec{g}_2) - Sim(b, \vec{g}_1) \right)$$

Sets of
Sentences

Prompt Engineering (Hard-Prompting)

I → ... Embedding → FF DNN → Output

Postprocess

Preprocess Intervention
Embedders
Learning
Pretraining Set

In Process Intervention
output
Fine-Tuning

---

Preprocess interventions
— Adding Sentences (CDA: Counterfactual Data Aug.)
   [He] is a doctor
      ↳ [She] is a doctor...

— Rewriting Sentences
   "[Name] follows his dream of being a doctor"
                    the
   "  "     "     "    "  "  ~  ~ "

— Removing Sentences: Removes the Problematic Sentences

Issues:

1) The Size of Pretraining Set is HUGE!
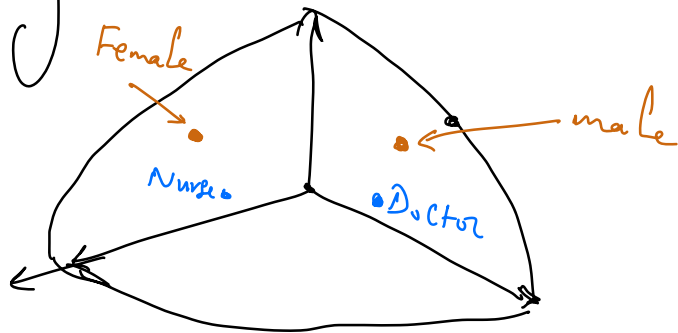   ↳ Debiasing the input for Fine Tuning ✓  ✗

2) Relying on Existing Tools is questionable!

3) Problematic for more than two groups
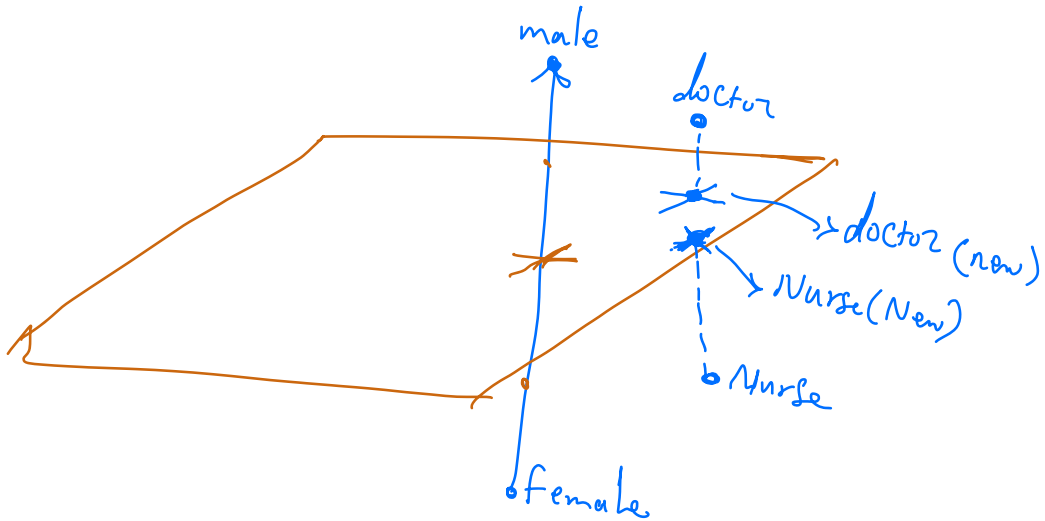   ↳ Significantly Change the data!

4) For unConsidered groups, the Entire learning Process Should repeat!

_____

Debiasing the Embeddings (InProcess Interventions)
   ↳ word Embedding



Goal: minimally change the Embeddings to remove the Bias

male

doctor
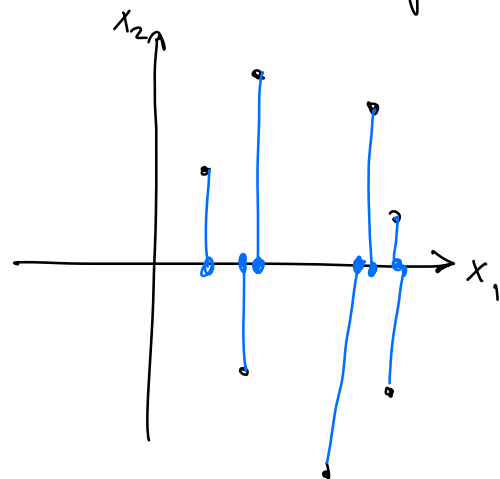
doctor (new)

Nurse (New)
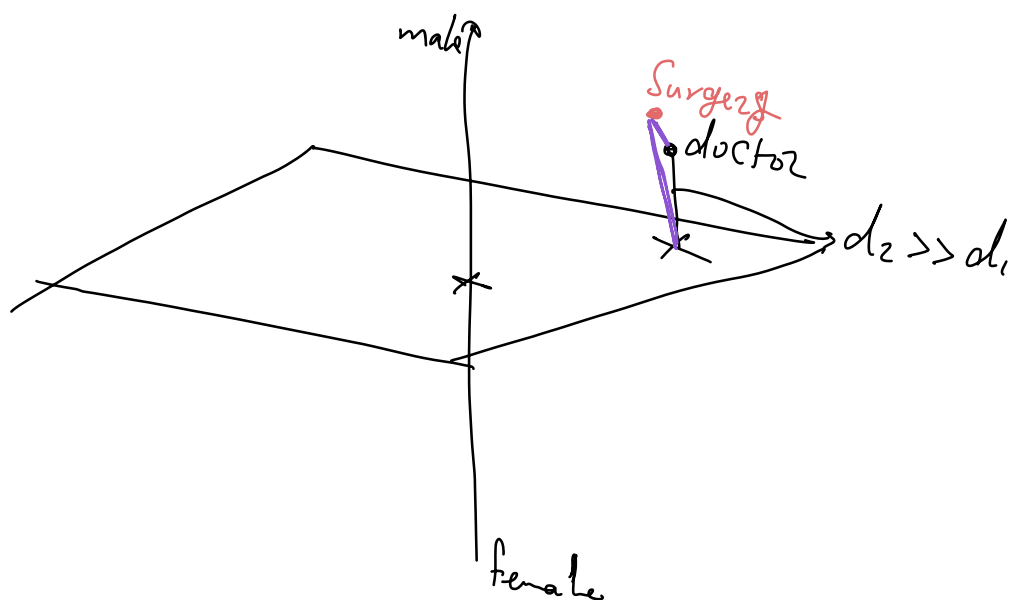
Nurse

female

Approach 1:  Project ALL the words.
↳) Reducing the dimension by 1
↳ missing some information
we lose useful associations with groups

Approach 2: Only Project the "Problematic" words
and don't change the others.
doctor, nurse,...



- we may mess up semantic similarities with projected words.