

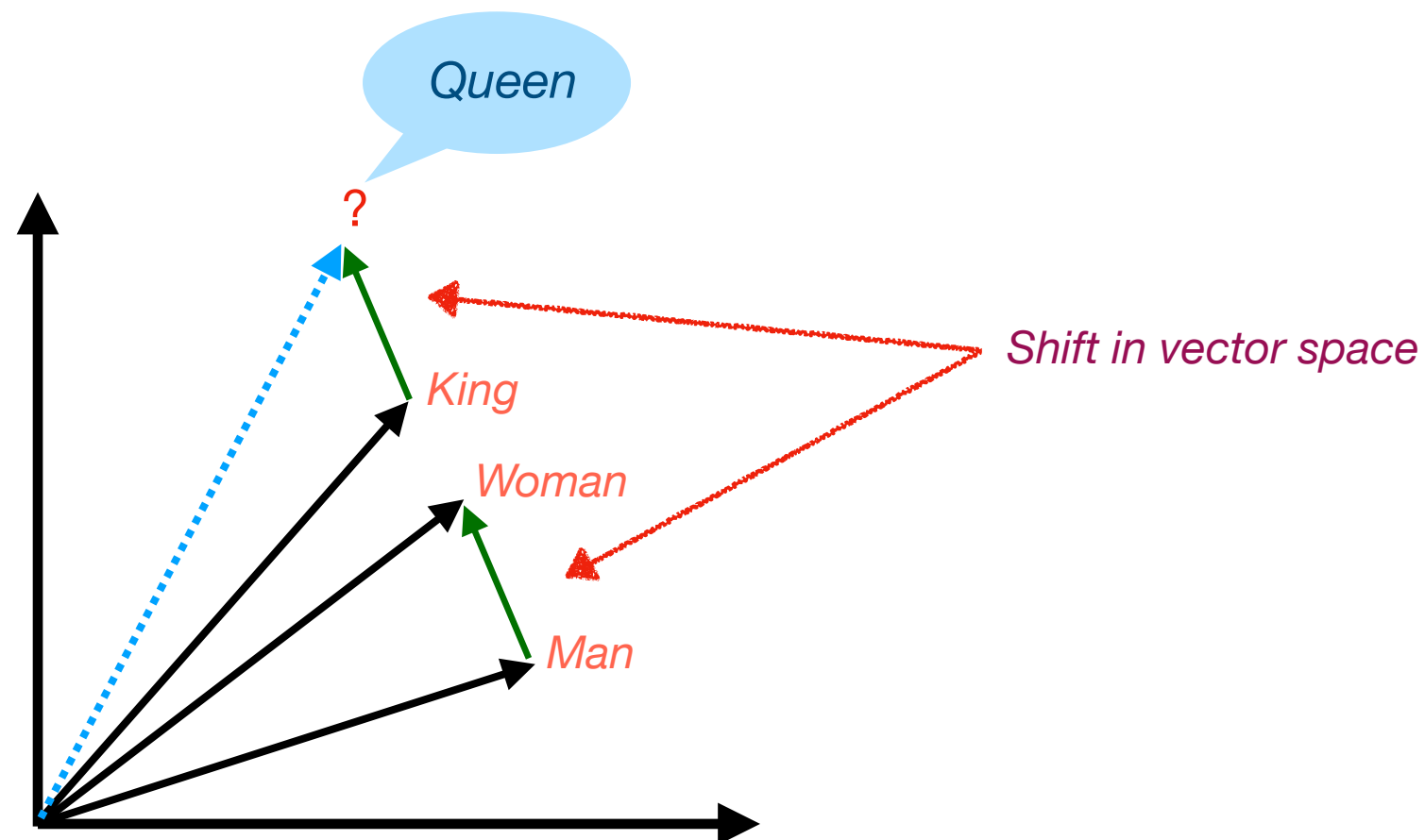
# Gender Bias in BERT

# Numeric Representations of Words

Word embedding: high-dimensional yet compact meaningful vector representations of words.

A meaningful vector representations of words can give answers to —

$$(\text{Woman} - \text{Man}) + \text{King} = ?$$



# Gender Bias in Word Representations

# Word-word co-occurrence

# Dataset



Male



Female

## Bias inherit

# Gender Bias in Word Representations

# Word-word co-occurrence

# Dataset

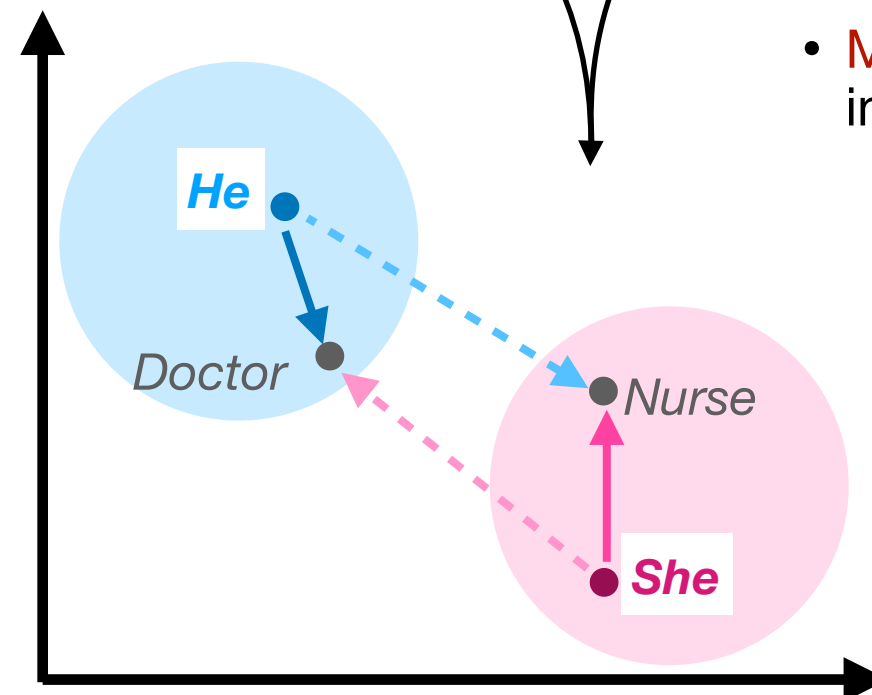


## Bias inherit

Male

Female

# Word vectors



- **Model** inherit bias from the dataset. For instance, unequal distances:

$$He - Doctor < She - Doctor$$

$$He - Nurse > She - Nurse$$

Do you notice gender **colour** **stereotype** here?

# Problem Formulation

Since it is hard to find and remove all the gender stereotypes from the dataset, we focus on identify biases inherited by the models trained on such datasets.

- Define: When do we call a system is gender biased?
- Quantify: How to measure/quantify the gender bias?
- Solve: How to gender debias?

# Measure Bias

## How to measure the bias?

- Compute the difference between similarities of a gender neutral word (eg., Doctor) from gender-specific words (eg. Male and Female).

Similarity score: Cosine angle, Euclidean distance.

- Compute the difference between prediction scores when a gender-specific word is switched to a word specifying some other gender.
- Some other way that you may figure out in future...

# Contextualised Word Embedding

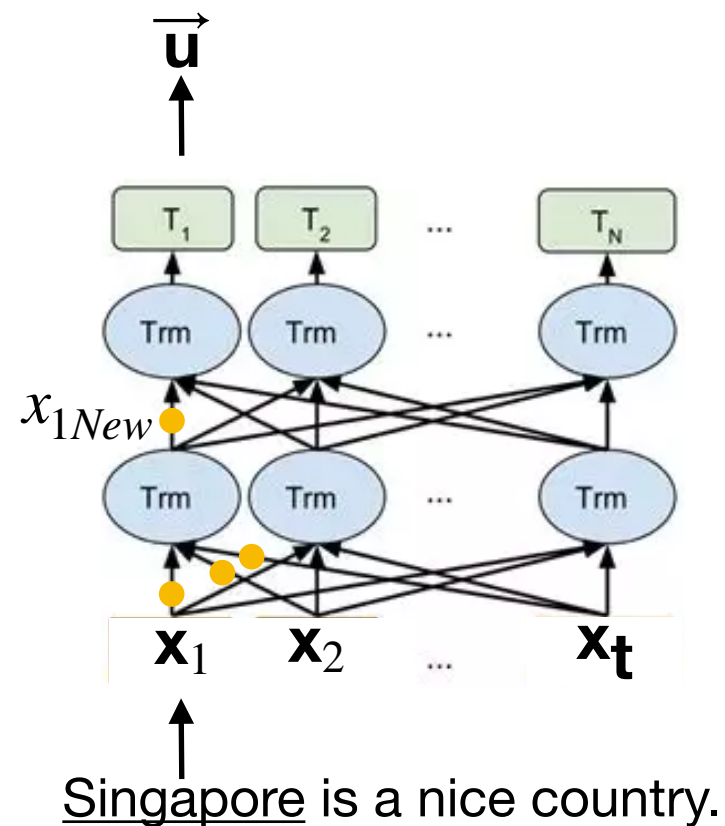
- BERT: Bidirectional Encoder Representations from Transformers (BERT) is popular provides contextualised word vectors. Consider two sentences:

S1: Singapore is a nice country.

S2: Singapore university of technology and design.

$\vec{u} = \text{BERT}(S1)$  (vector of Singapore in S1)

$\vec{v} = \text{BERT}(S2)$  (vector of Singapore in S2)





# Contextualised Word Embedding

- BERT: Bidirectional Encoder Representations from Transformers (BERT) is popular provides contextualised word vectors. Consider two sentences:

S1: Singapore is a nice country.

S2: Singapore university of technology and design.

$\vec{u} = \text{BERT}(\text{S1})$  (vector of Singapore in S1)

$\vec{v} = \text{BERT}(\text{S2})$  (vector of Singapore in S2)

We find...

- $\text{Cosine}(\vec{u}, \vec{v}) \neq 1$
- $\vec{u} - \vec{v} \neq 0$

Even though  $\vec{u}$  and  $\vec{v}$  are vectors for the same word Singapore, they are different because of the context.



# Measure Bias in Contextualised Embedding

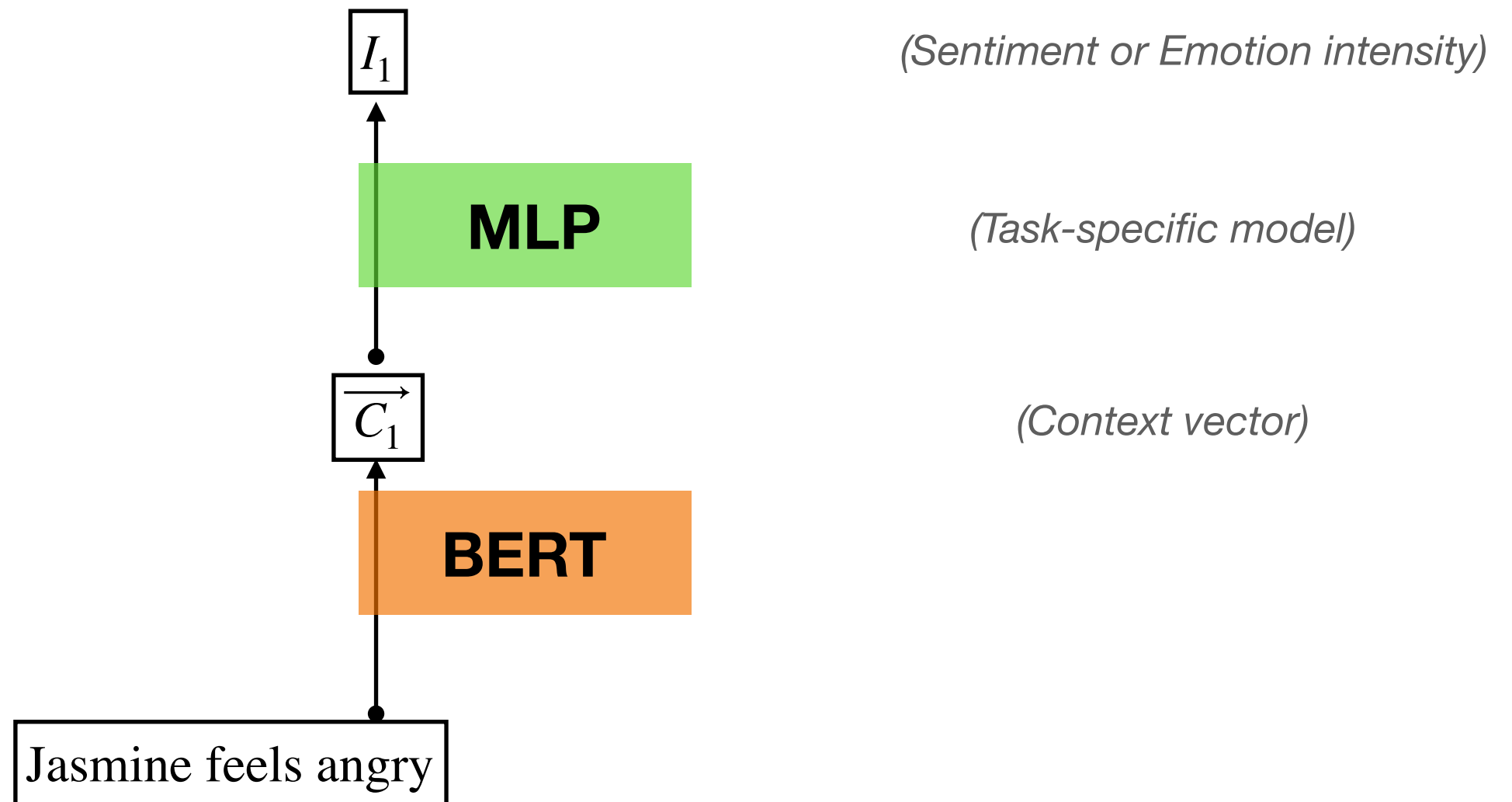
- One way to quantify gender bias in BERT is by observing predictions of downstream models utilising it as an underlying language model.

Downstream tasks we consider:

Given a sentence such as “Jasmine feels angry”, predict

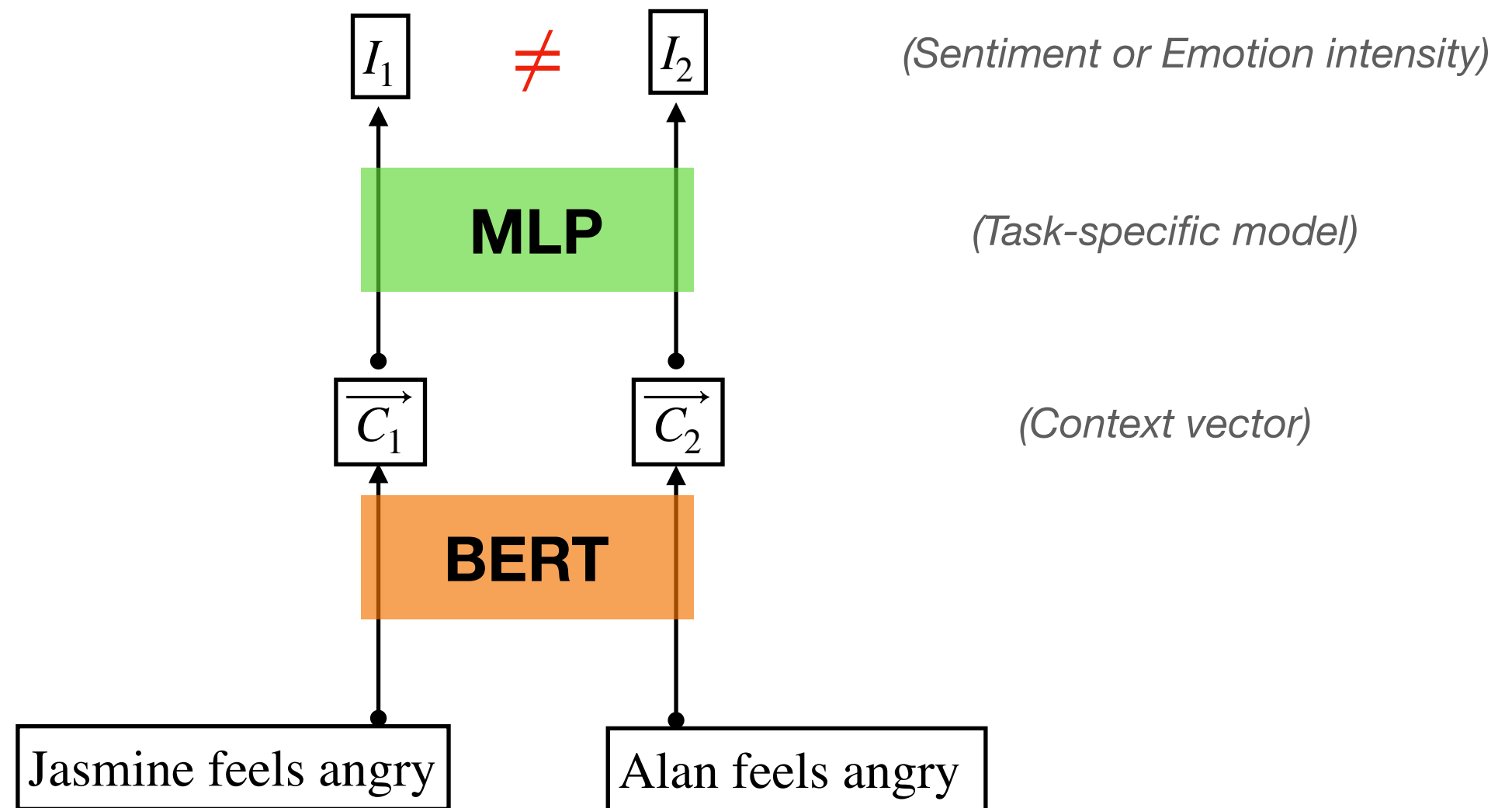
- 1) Intensity of emotion
- 2) Intensity of sentiment

# Measure Bias in Contextualised Embedding



Note: The MLP is kept simple (1-hidden layer) to prevent learning it's own bias.

# Measure Bias in Contextualised Embedding

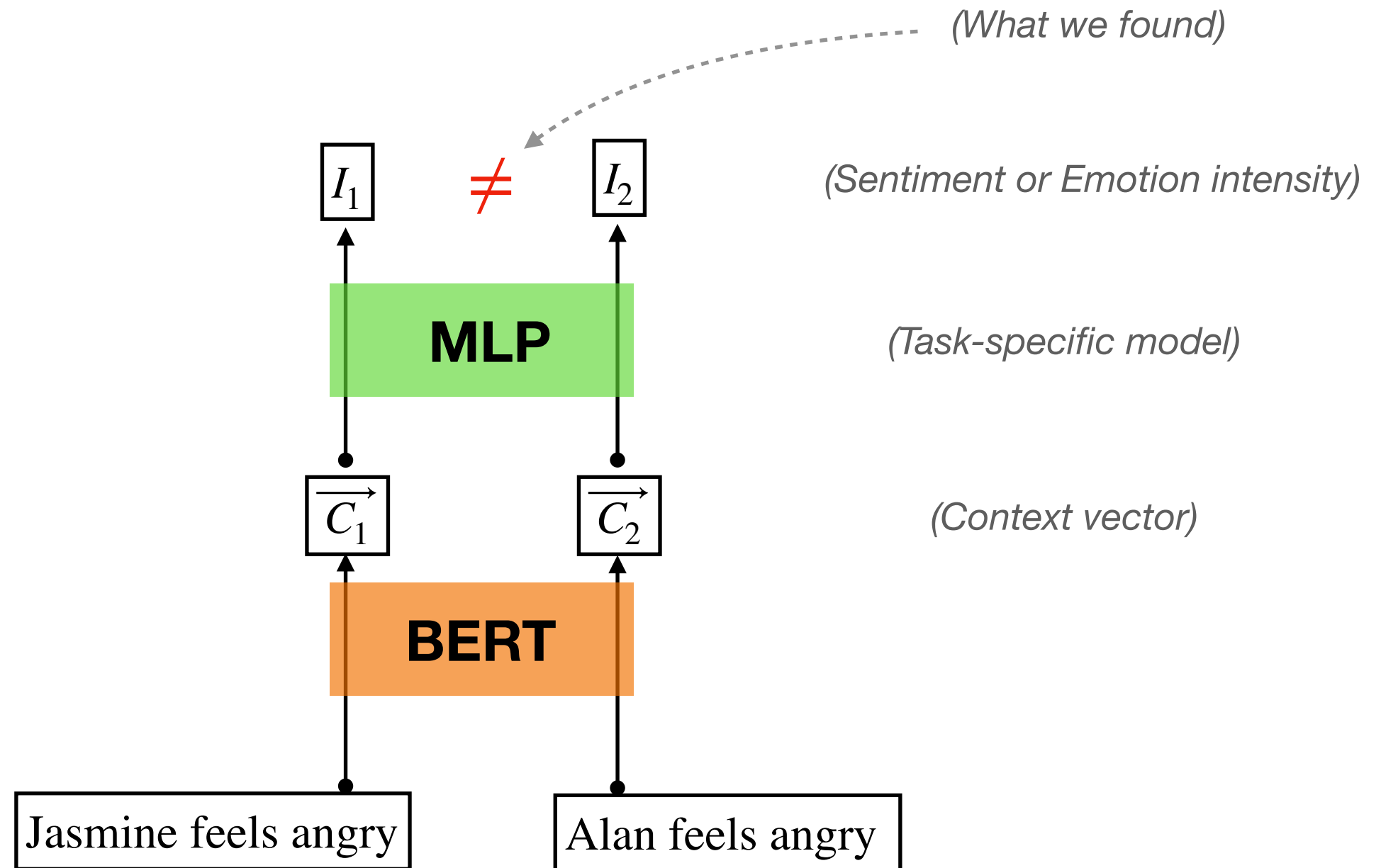


An ideal model is supposed to assign equal intensity scores.

**BERT-MLP does not follow that!**

Note: The MLP is kept simple (1-hidden layer) to prevent learning it's own bias.

# Measure Bias in Contextualised Embedding



An ideal model is supposed to assign equal intensity scores.

**BERT-MLP does not follow that!**

Note: The MLP is kept simple (1-hidden layer) to prevent learning it's own bias.

# Solve: Gender Debiasing Contextualised Embedding

We identify a subspace where BERT encodes gender information.

## Approach:

**Step1:** We collect a set of  $n$ -gender word pairs such as:

$$(f_i, m_i) \in \{(Queen, King), (She, He), \dots (Woman, Man)\}$$

*(The words in a gender pair are gender opposite of each other based on their common usage.)*

# Solve: Gender Debiasing Contextualised Embedding

We identify a subspace where BERT encodes gender information.

## Approach:

**Step1:** We collect a set of  $n$ -gender word pairs such as:

$$(f_i, m_i) \in \{(Queen, King), (She, He), \dots (Woman, Man)\}$$

*(The words in a gender pair are gender opposite of each other based on their common usage.)*

**Step2:** Feed them to BERT and obtain corresponding word embeddings.

$\vec{u}_i$  is word vector for  $f_i$  and  $\vec{v}_i$  is word vector for  $m_i$ .

# Solve: Gender Debiasing Contextualised Embedding

We identify a subspace where BERT encodes gender information.

## Approach:

**Step1:** We collect a set of  $n$ -gender word pairs such as:

$$(f_i, m_i) \in \{(Queen, King), (She, He), \dots (Woman, Man)\}$$

*(The words in a gender pair are gender opposite of each other based on their common usage.)*

**Step2:** Feed them to BERT and obtain corresponding word embeddings.

$\vec{u}_i$  is word vector for  $f_i$  and  $\vec{v}_i$  is word vector for  $m_i$ .

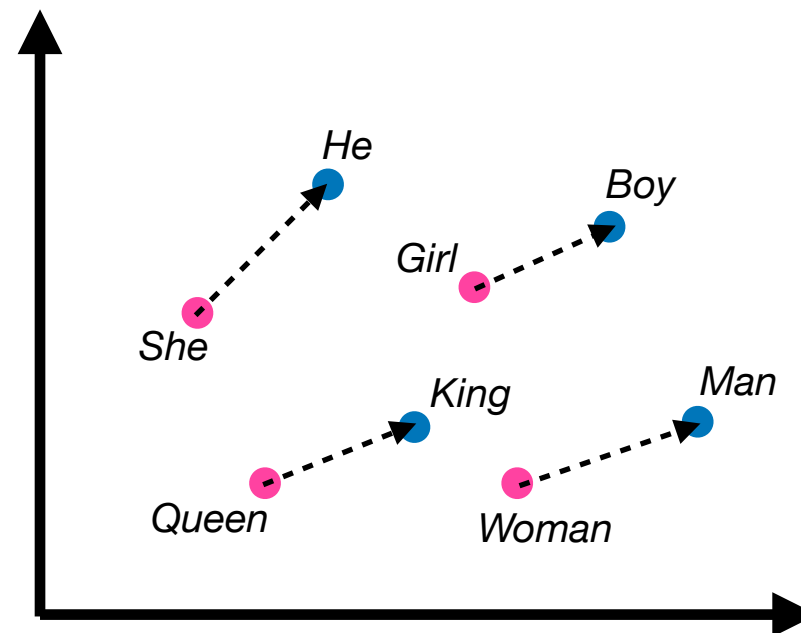
**Step3:** Obtain a set of difference vectors  $\{\vec{d}_1, \dots, \vec{d}_n\}$ , where  $\vec{d}_i = (\vec{v}_i - \vec{u}_i)$

*(Each difference vector shows shift from female to male in vector space.)*



# Solve: Gender Debiasing Contextualised Embedding

**Step4:** We have a set of  $n$ -difference vectors  $\{\vec{d}_1, \dots, \vec{d}_n\}$ , each vector shows a shift in gender direction.

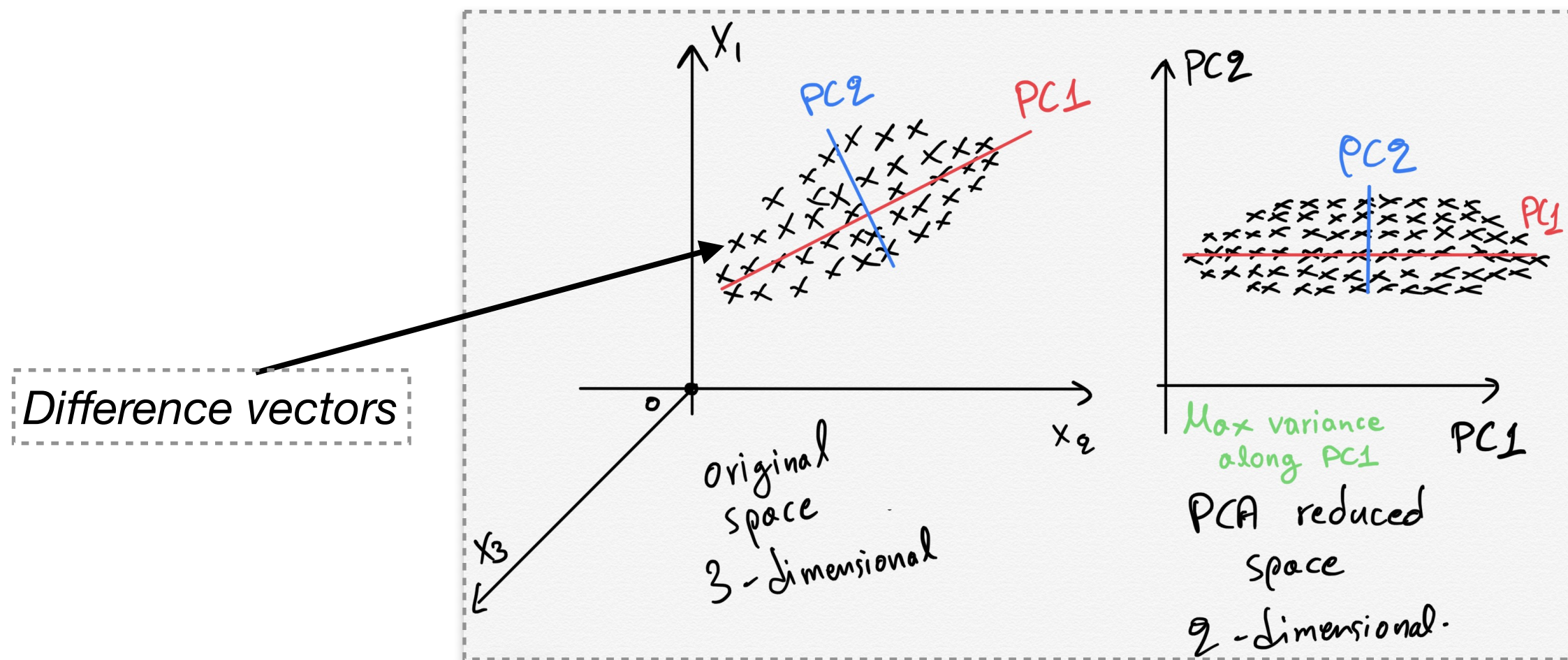


**Questions:** **Can** we use it to find the directions (subspace) that encode gender information? i.e., do we know any algorithm that can help uncover the principal directions which explain the most spread of vectors in  $\{\vec{d}_1, \dots, \vec{d}_n\}$ ?

# PCA Recap

Remember: First principal component 1<sub>st</sub> PC corresponds to the direction that best fits the data, 2<sub>nd</sub> PC is best fit direction orthogonal to 1<sub>st</sub> PC.

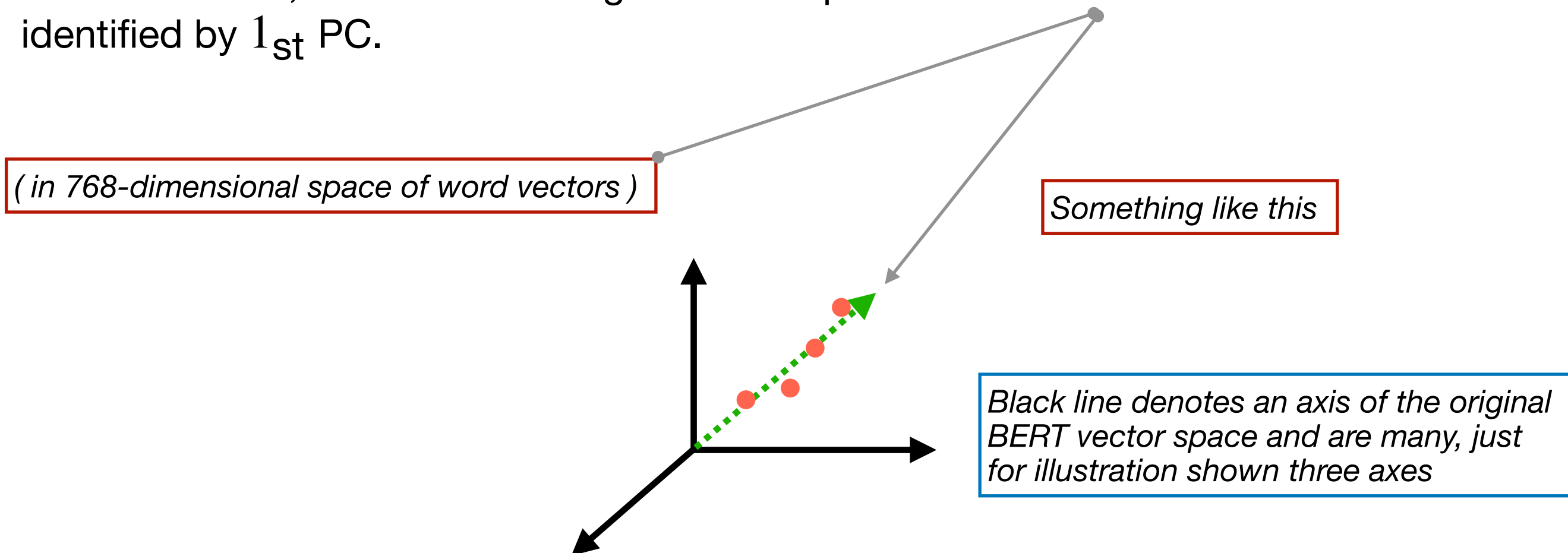
(Best fit implies the direction of a line that minimises the average square distance from data points to the line)



*Image Source*

# Solve: Gender Debiasing Contextualised Embedding

**Step5:** Compute PCs. The 1<sub>st</sub> PC covers explains most variance in the difference vectors. Hence, we consider the gender subspace to be 1-dimensional with axis identified by 1<sub>st</sub> PC.



**Step6:** Removing obtained gender direction during inference:

From each word vector produced by BERT, we remove the component in PC1. This will reduce the gender-specific information from the BERT word embedding.

---

**ALGORITHM 1:** Extracting layer-wise principal component in Gender subspace.

---

**Input** : - Strings pair  $(S_m, S_f)$ , which differ only in gender-specific words.

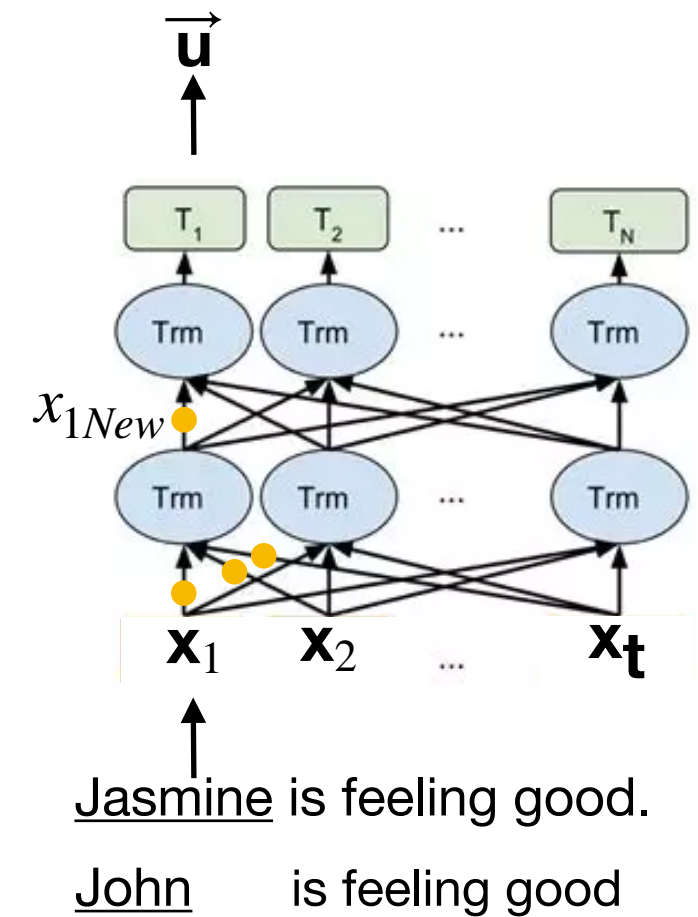
**Output** : -  $P$ =Layer-wise principal component set  $\{P_0, \dots, P_{12}\}$ .

```

1  $W_{tf} \leftarrow \text{Tokenize}(S_f)$            /* WP Tokenization */
2  $W_{tm} \leftarrow \text{Tokenize}(S_m)$        /* WP Tokenization */
3  $u_0 \leftarrow \text{Layer}_0(W_{tf})$          /* Context-independent input
   vectors for  $S_f$  */
4  $v_0 \leftarrow \text{Layer}_0(W_{tm})$          /* Context-independent input
   vectors for  $S_m$  */
5  $D_0 \leftarrow (v_0 - u_0)$              /* Difference vector */
6  $P_0 \leftarrow \text{PCA}(D_0)$              /* PC with maximum EV */
7 for  $j \leftarrow [1, 2, \dots, 12]$  do
8    $u_{j-1}^* \leftarrow \text{Proj}_{\perp P_{j-1}}(u_{j-1})$  /* Perpendicular
   projection */
9    $v_{j-1}^* \leftarrow \text{Proj}_{\perp P_{j-1}}(v_{j-1})$ 
10   $u_j \leftarrow \text{Layer}_j(u_{j-1}^*)$ 
11   $v_j \leftarrow \text{Layer}_j(v_{j-1}^*)$ 
12   $D_j \leftarrow (v_j - u_j)$            /* Difference vector */
13   $P_j \leftarrow \text{PCA}(D_j)$             $x_{1New}$ 
14 end

```

---



Algorithm: Iterative computation of principal components.

Emotion	Emotion Intensity						Valence Intensity					
	BERT			BERT <sup>De</sup>			BERT			BERT <sup>De</sup>		
	Pearson	$\Delta_{F\uparrow-M\downarrow}$	$\Delta_{M\uparrow-F\downarrow}$	Pearson	$\Delta_{F\uparrow-M\downarrow}(\%d)$	$\Delta_{M\uparrow-F\downarrow}(\%d)$	Pearson	$\Delta_{F\uparrow-M\downarrow}$	$\Delta_{M\uparrow-F\downarrow}$	Pearson	$\Delta_{F\uparrow-M\downarrow}(\%d)$	$\Delta_{M\uparrow-F\downarrow}(\%d)$
Joy	0.666	0.0396	0.0402	0.660	0.0143( $\downarrow$ <b>63.9</b> )	0.0143( $\downarrow$ <b>64.4</b> )	0.659	0.0346	0.0376	0.670	0.0209( $\downarrow$ <b>39.5</b> )	0.0138( $\downarrow$ <b>63.3</b> )
Fear	0.581	0.0202	0.0244	0.593	0.0152( $\downarrow$ <b>24.7</b> )	0.0158( $\downarrow$ <b>35.2</b> )		0.0263	0.0244		0.0156( $\downarrow$ <b>40.6</b> )	0.0123( $\downarrow$ <b>49.5</b> )
Sadness	0.615	0.0380	0.0138	0.604	0.0178( $\downarrow$ <b>58.9</b> )	0.0097( $\downarrow$ <b>29.7</b> )		0.0272	0.0205		0.0153( $\downarrow$ <b>43.7</b> )	0.0118( $\downarrow$ <b>42.4</b> )
Anger	0.627	0.0074	0.0316	0.626	0.0121( $\uparrow$ 63.5)	0.0149( $\downarrow$ <b>52.8</b> )		0.0219	0.0198		0.0130( $\downarrow$ <b>40.6</b> )	0.0119( $\downarrow$ <b>39.8</b> )

Table 2: Final-layer ( $layer_{12}$ ) of BERT and BERT<sup>De</sup> equity evaluation of the five-intensity regression models. %d refers to the percentage change in  $\Delta$  values. The p-values for 1) Emotion intensity models:  $\{\text{anger}\} \leq 0.05$  ( $\geq 0.70^*$ );  $\{\text{joy, fear, sad}\} \leq 0.20$  ( $\geq 0.70^*$ ). 2) The valence intensity model (emotion-wise p-values):  $\{\text{anger}\} \leq 0.05$  ( $\geq 0.75^*$ );  $\{\text{joy, fear, sad}\} \leq 0.20$  ( $\geq 0.70^*$ ), where values with \* denotes BERT<sup>De</sup>-based MLP regressor.

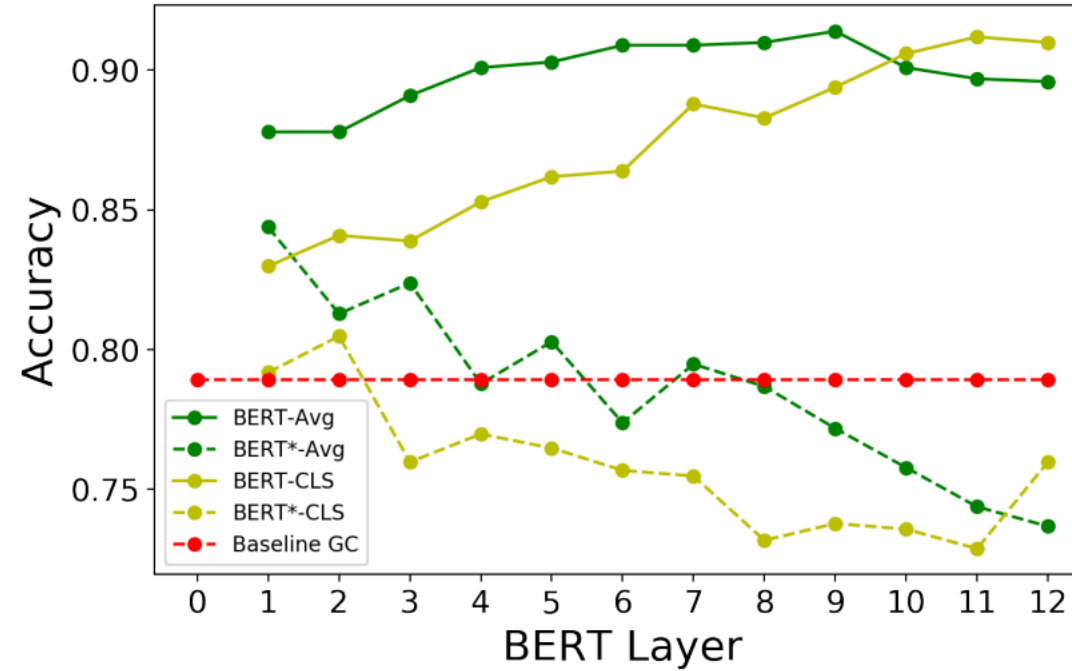


Figure 4: BERT-CLS and BERT<sup>De</sup>-CLS denote MLP accuracy using  $layer_k$  (x-axis) vectors in  $I_1$  setting. Similarly, BERT-Avg and BERT<sup>De</sup>-Avg refer to the  $I_2$  setting. Switching from BERT to BERT<sup>De</sup>, we see a significant drop in MLP gender-classification performance in both  $I_1$  and  $I_2$  inputs cases.