



# Analogies Explained

## Towards Understanding Word Embeddings

---

Carl Allen, Tim Hospedales

June 13 2019

School of Informatics, University of Edinburgh

# The Problem: linking semantics to geometry

from:

“man is to king as woman is to queen”

# The Problem: linking semantics to geometry

from:

“man is to king as woman is to queen”

explain:

$$\mathbf{w}_{king} - \mathbf{w}_{man} + \mathbf{w}_{woman} \approx \mathbf{w}_{queen}$$

# The Problem: linking semantics to geometry

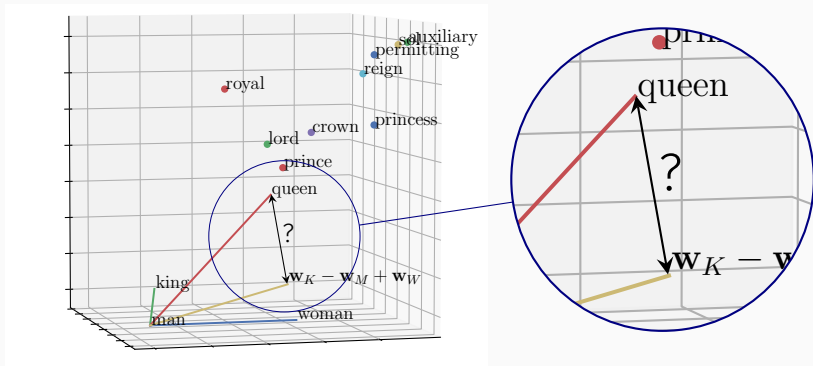
from:

“man is to king as woman is to queen”

explain:

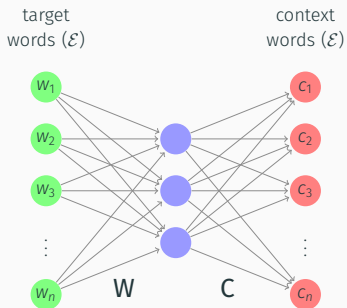
$$\mathbf{w}_{king} - \mathbf{w}_{man} + \mathbf{w}_{woman} \approx \mathbf{w}_{queen}$$

or rather:



# Word2Vec: SkipGram with Negative Sampling

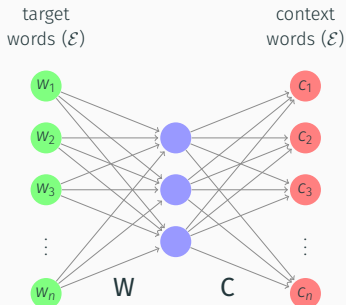
Mikolov et al. (2013a,b)



- $p(c_j|w_i)$  by **softmax** expensive

# Word2Vec: SkipGram with Negative Sampling

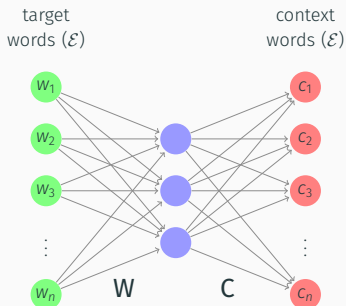
Mikolov et al. (2013a,b)



- $p(c_j|w_i)$  by **softmax** expensive
- use **sigmoid** with negative sampling ( $k$ )

# Word2Vec: SkipGram with Negative Sampling

Mikolov et al. (2013a,b)

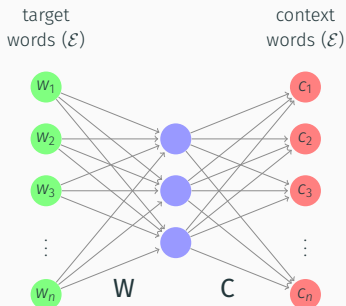


- $p(c_j|w_i)$  by **softmax** expensive
- use **sigmoid** with negative sampling ( $k$ )
- Levy and Goldberg (2014)

$$\begin{aligned} \mathbf{w}_i^T \mathbf{c}_j &\approx \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)} - \log k \\ &= \text{PMI}(w_i, c_j) - \log k \end{aligned}$$

# Word2Vec: SkipGram with Negative Sampling

Mikolov et al. (2013a,b)



- $p(c_j|w_i)$  by **softmax** expensive
- use **sigmoid** with negative sampling ( $k$ )
- Levy and Goldberg (2014)

$$\begin{aligned} \mathbf{w}_i^\top \mathbf{c}_j &\approx \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)} - \log k \\ &= \text{PMI}(w_i, c_j) - \log k \end{aligned}$$

$$\mathbf{W}^\top \mathbf{C} \approx \text{PMI} - \log k$$



“man is to king as woman is to queen”

“man is to king as woman is to queen”



man transforms to king as woman transforms to queen

“man is to king as woman is to queen”



man transforms to king as woman transforms to queen



{woman, king} paraphrases {man, queen}

“man is to king as woman is to queen”



man transforms to king as woman transforms to queen



{woman, king} paraphrases {man, queen}



$$PMI_{king} - PMI_{man} + PMI_{woman} \approx PMI_{queen}$$

“man is to king as woman is to queen”



man transforms to king as woman transforms to queen



{woman, king} paraphrases {man, queen}



$$PMI_{king} - PMI_{man} + PMI_{woman} \approx PMI_{queen}$$



$$W_{king} - W_{man} + W_{woman} \approx W_{queen}$$

“man is to king as woman is to queen”



man transforms to king as woman transforms to queen



{woman, king} paraphrases {man, queen}



$$\text{PMI}_{king} - \text{PMI}_{man} + \text{PMI}_{woman} \approx \text{PMI}_{queen}$$



$$\mathbf{W}_{king} - \mathbf{W}_{man} + \mathbf{W}_{woman} \approx \mathbf{W}_{queen}$$

“man is to king as woman is to queen”



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}



$$\text{PMI}_{king} - \text{PMI}_{man} + \text{PMI}_{woman} \approx \text{PMI}_{queen}$$



$$\mathbf{w}_{king} - \mathbf{w}_{man} + \mathbf{w}_{woman} \approx \mathbf{w}_{queen}$$

“man is to king as woman is to queen”



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}



$$\text{PMI}_{king} - \text{PMI}_{man} + \text{PMI}_{woman} \approx \text{PMI}_{queen}$$



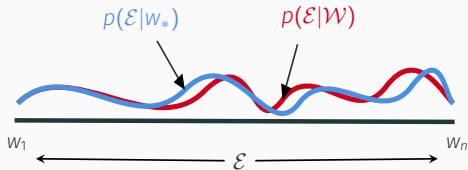
$$\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$$

$$\mathbf{w}_{king} - \mathbf{w}_{man} + \mathbf{w}_{woman} \approx \mathbf{w}_{queen}$$



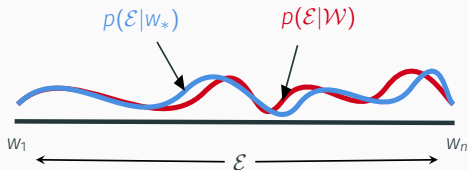
# Paraphrase<sup>†</sup> of $\mathcal{W}$ by $w_*$

Intuition: word  $w_* \in \mathcal{E}$  **paraphrases** word set  $\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$ , if  $w_*$  and  $\mathcal{W}$  are *semantically interchangeable*.



## Paraphrase<sup>†</sup> of $\mathcal{W}$ by $w_*$

Intuition: word  $w_* \in \mathcal{E}$  **paraphrases** word set  $\mathcal{W} = \{w_1, \dots, w_m\} \subseteq \mathcal{E}$ , if  $w_*$  and  $\mathcal{W}$  are *semantically interchangeable*.



**Definition (D1):**  $w_* \in \mathcal{E}$  paraphrases  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}| < l$ , if paraphrase error  $\rho^{w, w_*} \in \mathbb{R}^n$  is (element-wise) small:

$$\rho_j^{w, w_*} = \log \frac{p(c_j|w_*)}{p(c_j|\mathcal{W})}, \quad c_j \in \mathcal{E}$$

<sup>†</sup>Inspired by Gittens et al. (2017)

## Summing PMI vectors of a paraphrase

$$\text{PMI}_1 + \text{PMI}_2 \approx \text{PMI}_* ?$$

## Summing PMI vectors of a paraphrase

$\text{PMI}_1 + \text{PMI}_2 \approx \text{PMI}_* ?$

$$\text{PMI}(w_*, c_j) - (\text{PMI}(w_1, c_j) + \text{PMI}(w_2, c_j))$$

## Summing PMI vectors of a paraphrase

$\text{PMI}_1 + \text{PMI}_2 \approx \text{PMI}_* ?$

$$\begin{aligned} & \text{PMI}(w_*, c_j) - (\text{PMI}(w_1, c_j) + \text{PMI}(w_2, c_j)) \\ &= \log \frac{p(w_* | c_j)}{p(w_*)} - \log \frac{p(w_1 | c_j)p(w_2 | c_j)}{p(w_1)p(w_2)} \end{aligned}$$

## Summing PMI vectors of a paraphrase

$\text{PMI}_1 + \text{PMI}_2 \approx \text{PMI}_* ?$

$$\begin{aligned} & \text{PMI}(w_*, c_j) - (\text{PMI}(w_1, c_j) + \text{PMI}(w_2, c_j)) \\ &= \log \frac{p(w_*|c_j)}{p(w_*)} - \log \frac{p(w_1|c_j)p(w_2|c_j)}{p(w_1)p(w_2)} + \log \frac{p(\mathcal{W}|c_j)}{p(\mathcal{W})} + \log \frac{p(\mathcal{W})}{p(\mathcal{W})} \end{aligned}$$

# Summing PMI vectors of a paraphrase

$$\text{PMI}_1 + \text{PMI}_2 \approx \text{PMI}_* ?$$

$$\text{PMI}(w_*, c_j) - (\text{PMI}(w_1, c_j) + \text{PMI}(w_2, c_j))$$

$$= \log \frac{p(w_*|c_j)}{p(w_*)} - \log \frac{p(w_1|c_j)p(w_2|c_j)}{p(w_1)p(w_2)} + \log \frac{p(\mathcal{W}|c_j)}{p(\mathcal{W})} + \log \frac{p(\mathcal{W})}{p(\mathcal{W})}$$

$$= \underbrace{\log \frac{p(c_j|w_*)}{p(c_j|\mathcal{W})}}_{\rho_j^{\mathcal{W}, w_*}} + \underbrace{\log \frac{p(\mathcal{W}|c_j)}{p(w_1|c_j)p(w_2|c_j)}}_{\sigma_j^{\mathcal{W}}} - \underbrace{\log \frac{p(\mathcal{W})}{p(w_1)p(w_2)}}_{\tau^{\mathcal{W}}}$$

paraphrase  
error

conditional  
independence  
error

independence  
error

# Summing PMI vectors of a paraphrase

$\text{PMI}_1 + \text{PMI}_2 \approx \text{PMI}_* ?$

$$\begin{aligned} & \text{PMI}(w_*, c_j) - (\text{PMI}(w_1, c_j) + \text{PMI}(w_2, c_j)) \\ &= \log \frac{p(w_*|c_j)}{p(w_*)} - \log \frac{p(w_1|c_j)p(w_2|c_j)}{p(w_1)p(w_2)} + \log \frac{p(\mathcal{W}|c_j)}{p(\mathcal{W}|c_j)} + \log \frac{p(\mathcal{W})}{p(\mathcal{W})} \\ &= \underbrace{\log \frac{p(c_j|w_*)}{p(c_j|\mathcal{W})}}_{\rho_j^{\mathcal{W}, w_*}} + \underbrace{\log \frac{p(\mathcal{W}|c_j)}{p(w_1|c_j)p(w_2|c_j)}}_{\sigma_j^{\mathcal{W}}} - \underbrace{\log \frac{p(\mathcal{W})}{p(w_1)p(w_2)}}_{\tau^{\mathcal{W}}} \\ & \qquad \text{paraphrase error} \qquad \qquad \text{conditional independence error} \qquad \qquad \text{independence error} \end{aligned}$$

**Lemma 1:** For any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}| < l$ :

$$\text{PMI}_* = \sum_{w \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1}$$



## Generalised Paraphrase (of $\mathcal{W}$ by $\mathcal{W}_*$ )

**Lemma 1:** For any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}| < l$ :

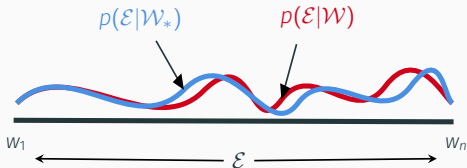
$$\text{PMI}_* = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1}$$

# Generalised Paraphrase (of $\mathcal{W}$ by $\mathcal{W}_*$ )

**Lemma 1:** For any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}| < l$ :

$$\text{PMI}_* = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{\mathcal{W}, w_*} + \sigma^{\mathcal{W}} - \tau^{\mathcal{W}} \mathbf{1}$$

Replace word  $w_*$  with word set  $\mathcal{W}_* \subseteq \mathcal{E}$ :

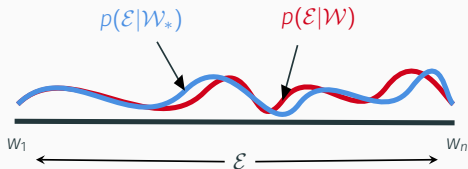


# Generalised Paraphrase (of $\mathcal{W}$ by $\mathcal{W}_*$ )

**Lemma 1:** For any word  $w_* \in \mathcal{E}$  and word set  $\mathcal{W} \subseteq \mathcal{E}$ ,  $|\mathcal{W}| < l$ :

$$\text{PMI}_* = \sum_{w \in \mathcal{W}} \text{PMI}_i + \rho^{w, w_*} + \sigma^w - \tau^{w*} \mathbf{1}$$

Replace word  $w_*$  with word set  $\mathcal{W}_* \subseteq \mathcal{E}$ :



**Lemma 2:** For any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$\sum_{w \in \mathcal{W}_*} \text{PMI}_i = \sum_{w \in \mathcal{W}} \text{PMI}_i + \rho^{w, w_*} + \sigma^w - \sigma^{w_*} - (\tau^w - \tau^{w_*}) \mathbf{1}$$

## Paraphrase: the link from semantics to geometry

Lemma 2: For any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$\sum_{w \in \mathcal{W}_*} \text{PMI}_i = \sum_{w \in \mathcal{W}} \text{PMI}_i + \rho^{w, w_*} + \sigma^w - \sigma^{w_*} - (\tau^w - \tau^{w_*}) \mathbf{1}$$

## Paraphrase: the link from semantics to geometry

Lemma 2: For any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$\sum_{w_i \in \mathcal{W}_*} \text{PMI}_i = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{w, w_*} + \sigma^w - \sigma^{w_*} - (\tau^w - \tau^{w_*}) \mathbf{1}$$

So, if

$\mathcal{W} = \{\text{woman, king}\}$  paraphrases  $\mathcal{W}_* = \{\text{man, queen}\}$ ,

# Paraphrase: the link from semantics to geometry

Lemma 2: For any word sets  $\mathcal{W}, \mathcal{W}_* \subseteq \mathcal{E}$ ,  $|\mathcal{W}|, |\mathcal{W}_*| < l$ :

$$\sum_{w_i \in \mathcal{W}_*} \text{PMI}_i = \sum_{w_i \in \mathcal{W}} \text{PMI}_i + \rho^{w, w_*} + \sigma^w - \sigma^{w_*} - (\tau^w - \tau^{w_*})\mathbf{1}$$

So, if

$\mathcal{W} = \{\text{woman}, \text{king}\}$  paraphrases  $\mathcal{W}_* = \{\text{man}, \text{queen}\}$ ,

then:

$$\text{PMI}_{\text{queen}} \approx \text{PMI}_{\text{king}} - \text{PMI}_{\text{man}} + \text{PMI}_{\text{woman}} + \underbrace{\sigma^w - \sigma^{w_*} - (\tau^w - \tau^{w_*})\mathbf{1}}_{\text{net dependence error}}$$

“man is to king as woman is to queen”



man *transforms to* king as woman *transforms to* queen



{woman, king} *paraphrases* {man, queen}

↓ dependence error

$$PMI_{king} - PMI_{man} + PMI_{woman} \approx PMI_{queen}$$

↓  $PMI_i \approx \mathbf{w}_i^T \mathbf{C}$

$$\mathbf{W}_{king} - \mathbf{W}_{man} + \mathbf{W}_{woman} \approx \mathbf{W}_{queen}$$

## Word Transformation: a change of perspective

A paraphrase  $w_*$  of  $\mathcal{W}$  can be thought of as a **word transformation** from some  $w \in \mathcal{W}$  to  $w_*$  by *adding*  $\mathcal{W}^+ = \{w_i \in \mathcal{W}, w_i \neq w\}$ , e.g.

$$\{man, royal\} \approx_p king \quad \Longrightarrow \quad man \xrightarrow{+royal} king$$



## Word Transformation: a change of perspective

A paraphrase  $w_*$  of  $\mathcal{W}$  can be thought of as a **word transformation** from some  $w \in \mathcal{W}$  to  $w_*$  by *adding*  $\mathcal{W}^+ = \{w_i \in \mathcal{W}, w_i \neq w\}$ , e.g.

$$\{man, royal\} \approx_p king \quad \Longrightarrow \quad man \xrightarrow{+royal} king$$

Added words **contextualise**  $w$ , such that the induced distribution better aligns with that of  $w_*$ .

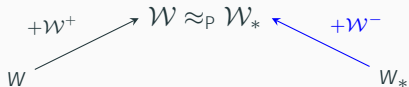
# Word Transformation: a change of perspective

A paraphrase  $w_*$  of  $\mathcal{W}$  can be thought of as a **word transformation** from some  $w \in \mathcal{W}$  to  $w_*$  by *adding*  $\mathcal{W}^+ = \{w_i \in \mathcal{W}, w_i \neq w\}$ , e.g.

$$\{man, royal\} \approx_p king \quad \Longrightarrow \quad man \xrightarrow{+royal} king$$

Added words **contextualise**  $w$ , such that the induced distribution better aligns with that of  $w_*$ .

Paraphrase  $\mathcal{W}$  by  $\mathcal{W}_*$  can be thought of as a word transformation from some  $w \in \mathcal{W}$  to some  $w_* \in \mathcal{W}_*$  by adding to both ...



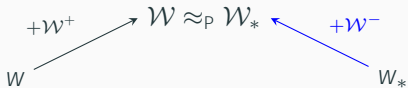
# Word Transformation: a change of perspective

A paraphrase  $w_*$  of  $\mathcal{W}$  can be thought of as a **word transformation** from some  $w \in \mathcal{W}$  to  $w_*$  by *adding*  $\mathcal{W}^+ = \{w_i \in \mathcal{W}, w_i \neq w\}$ , e.g.

$$\{man, royal\} \approx_p king \quad \Longrightarrow \quad man \xrightarrow{+royal} king$$

Added words **contextualise**  $w$ , such that the induced distribution better aligns with that of  $w_*$ .

Paraphrase  $\mathcal{W}$  by  $\mathcal{W}_*$  can be thought of as a word transformation from some  $w \in \mathcal{W}$  to some  $w_* \in \mathcal{W}_*$  by adding to both ...



... or adding to one side and *subtracting* from the other:



## Word Transformation: a change of perspective (cont.)



A generalised paraphrase *is* a word transformation from  $w \in \mathcal{W}$  to  $w_* \in \mathcal{W}_*$ , where:

- added words *narrow* context
- subtracted words *broaden* context

## Word Transformation: a change of perspective (cont.)



A generalised paraphrase *is* a word transformation from  $w \in \mathcal{W}$  to  $w_* \in \mathcal{W}_*$ , where:

- added words *narrow* context
- subtracted words *broaden* context

Providing a “richer dictionary” to explain the difference between  $w$  and  $w_*$ , or rather, how “**w is to w\_\***”.

## Word Transformation: a change of perspective (cont.)



A generalised paraphrase *is* a word transformation from  $w \in \mathcal{W}$  to  $w_* \in \mathcal{W}_*$ , where:

- added words *narrow* context
- subtracted words *broaden* context

Providing a “richer dictionary” to explain the difference between  $w$  and  $w_*$ , or rather, how “**w is to w\***”.

**Definition (D4):** We say “ $w_a$  is to  $w_{a^*}$  as  $w_b$  is to  $w_{b^*}$ ” iff there exist  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform  $w_a$  to  $w_{a^*}$  and  $w_b$  to  $w_{b^*}$ .

## Word Transformation: a change of perspective (cont.)

That is, we say:

“man is to king as woman is to queen”

*iff* there exist  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform man to king and woman to queen.

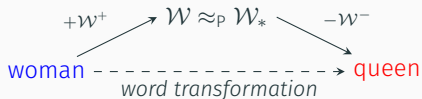
# Word Transformation: a change of perspective (cont.)

That is, we say:

“man is to king as woman is to queen”

iff there exist  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform man to king and woman to queen.

That is:





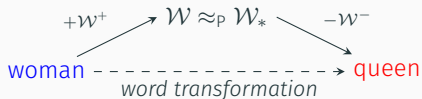
## Word Transformation: a change of perspective (cont.)

That is, we say:

“**man** is to **king** as **woman** is to **queen**”

iff there exist  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform **man** to **king** and **woman** to **queen**.

That is:



Let  $\mathcal{W}^+ = \{\text{king}\}$ ,  $\mathcal{W}^- = \{\text{man}\}$ .

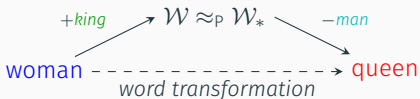
## Word Transformation: a change of perspective (cont.)

That is, we say:

“**man** is to **king** as **woman** is to **queen**”

iff there exist  $\mathcal{W}^+, \mathcal{W}^- \subseteq \mathcal{E}$  that simultaneously transform **man** to **king** and **woman** to **queen**.

That is:



Let  $\mathcal{W}^+ = \{\text{king}\}$ ,  $\mathcal{W}^- = \{\text{man}\}$ .

“man is to king as woman is to queen”



man transforms to king as woman transforms to queen



{woman, king} paraphrases {man, queen}

$\Downarrow$  dependence error

$$\text{PMI}_{king} - \text{PMI}_{man} + \text{PMI}_{woman} \approx \text{PMI}_{queen}$$

$\Downarrow$   $\text{PMI}_i \approx \mathbf{w}_i^T \mathbf{C}$

$$\mathbf{W}_{king} - \mathbf{W}_{man} + \mathbf{W}_{woman} \approx \mathbf{W}_{queen}$$

## The Solution: linking semantics to geometry

“man is to king as woman is to queen”

## The Solution: linking semantics to geometry

“man is to king as woman is to queen”

implies:

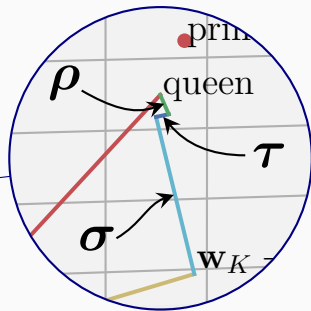
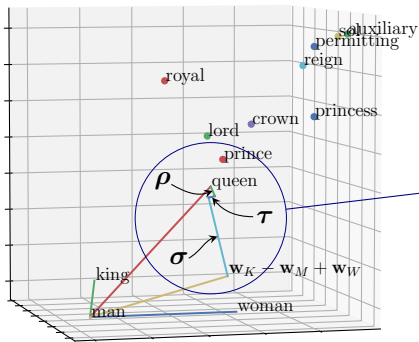
$$\mathbf{W}_{king} - \mathbf{W}_{man} + \mathbf{W}_{woman} \stackrel{\rho, \sigma, \tau}{\approx} \mathbf{W}_{queen}$$

# The Solution: linking semantics to geometry

“man is to king as woman is to queen”

implies:

$$\mathbf{W}_{king} - \mathbf{W}_{man} + \mathbf{W}_{woman} \stackrel{\rho, \sigma, \tau}{\approx} \mathbf{W}_{queen}$$



## References

---

- Alex Gittens, Dimitris Achlioptas, and Michael W Mahoney. Skip-gram-zipf+ uniform= vector additivity. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 69–76, 2017.
- Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *Advances in neural information processing systems*, 2014.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013a.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, 2013b.