

# Interpreting KG Relation Representation from Word Embeddings

Carl Allen<sup>1\*</sup> Ivana Balažević<sup>1\*</sup> Timothy Hospedales<sup>1,2</sup>

{carl.allen, ivana.balazevic, t.hospedales}@ed.ac.uk

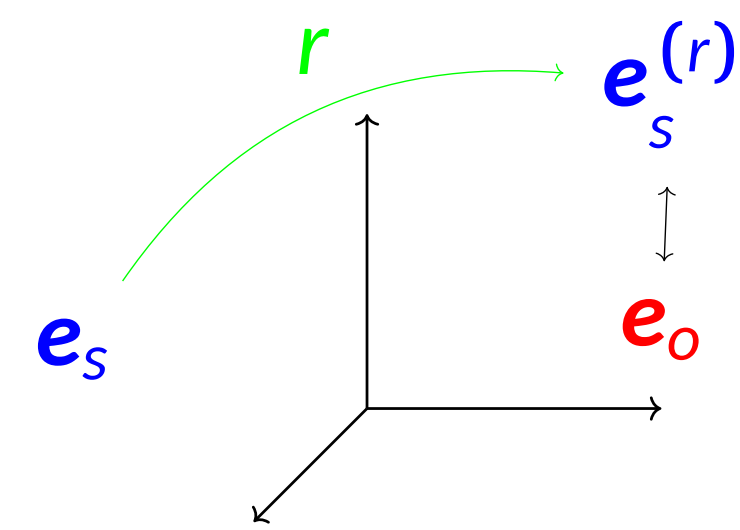
<sup>1</sup> School of Informatics, University of Edinburgh, UK <sup>2</sup> Samsung AI Centre, Cambridge, UK

## Abstract: How do KG Embeddings Capture Semantics?

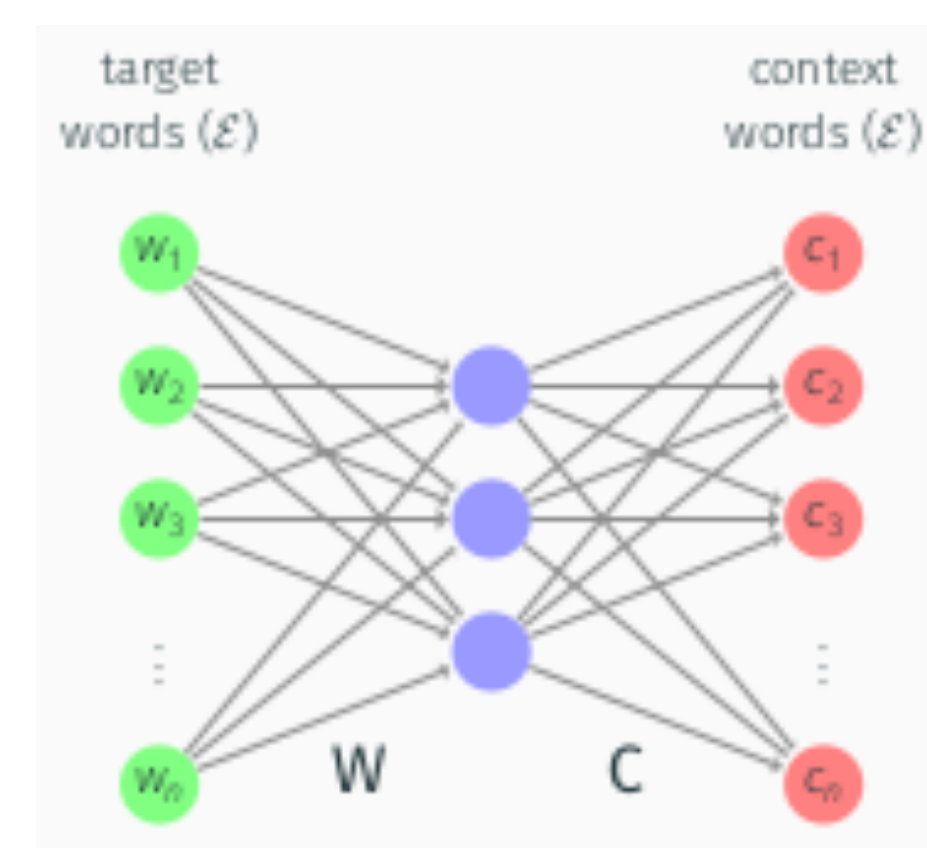
- There are many knowledge graph (KG) representation models, yet **little is understood of the latent structure they learn**.
- Building on recent **word embedding** theory (Allen and Hospedales, 2019; Allen et al., 2019), we derive **geometric properties** of knowledge graph relation representations (**relation conditions**) that map word embeddings of a subject entity to those of related object entities.
- We show that the better a model's architecture satisfies a relation's conditions, the better its performance at link prediction.

## Background: Knowledge Graph Representation

- KGs store facts: binary **relations** between **entities** ( $e_s, r, e_o$ ).
- Enable computational reasoning over KGs, e.g. question answering and inferring new facts (**link prediction**).
- Requires representation, typically:
  - each entity by a vector **embedding**  $e \in \mathbb{R}^d$ ,
  - each relation by a **transformation** from subject entity to object entity,
- A **score function** measures proximity between transformed entity embeddings.



## Simplify: Consider Word Embeddings



$\mathcal{E}$  = dictionary of all words

**Word2vec** (Mikolov et al., 2013) loss function is minimised when:

- **Low-rank case** (Levy and Goldberg, 2014):

$$w_i^T c_j = \log \frac{p(w_i, c_j)}{p(w_i)p(c_j)} - \log k \doteq s_{i,j}$$

- **General case** (Allen et al., 2019): Word embeddings  $w_i$  are a (non-linear) projection of **PMI vectors** ( $p^i$ ).

$$\text{PMI vectors: } p^j = \left\{ \log \frac{p(c_j | w_i)}{p(c_j)} \right\}_{c_j \in \mathcal{E}} = \log \frac{p(\mathcal{E} | w_i)}{p(\mathcal{E})}$$

## PMI Vector Interactions

**Similarity:** similar words induce similar distributions over context words. Subtraction of PMI vectors finds similarity (e.g. synonyms):

$$p^i - p^j = \log \frac{p(\mathcal{E} | w_i)}{p(\mathcal{E} | w_j)} = \rho^{i,j}$$

**Paraphrases:** word sets with similar aggregate semantic meaning, e.g. {man, royal}  $\approx$  king. Addition of PMI vectors finds paraphrases:

$$p^i + p^j = \log \frac{p(\mathcal{E} | w_i)}{p(\mathcal{E})} + \log \frac{p(\mathcal{E} | w_j)}{p(\mathcal{E})} \\ = p^k + \underbrace{\log \frac{p(\mathcal{E} | w_i, w_j)}{p(\mathcal{E} | w_k)}}_{\rho^{\{i,j\},k}} - \underbrace{\log \frac{p(w_i, w_j | \mathcal{E})}{p(w_i | \mathcal{E})p(w_j | \mathcal{E})}}_{\sigma^{i,j}} + \underbrace{\log \frac{p(w_i, w_j)}{p(w_i)p(w_j)}}_{\tau^{i,j}}$$

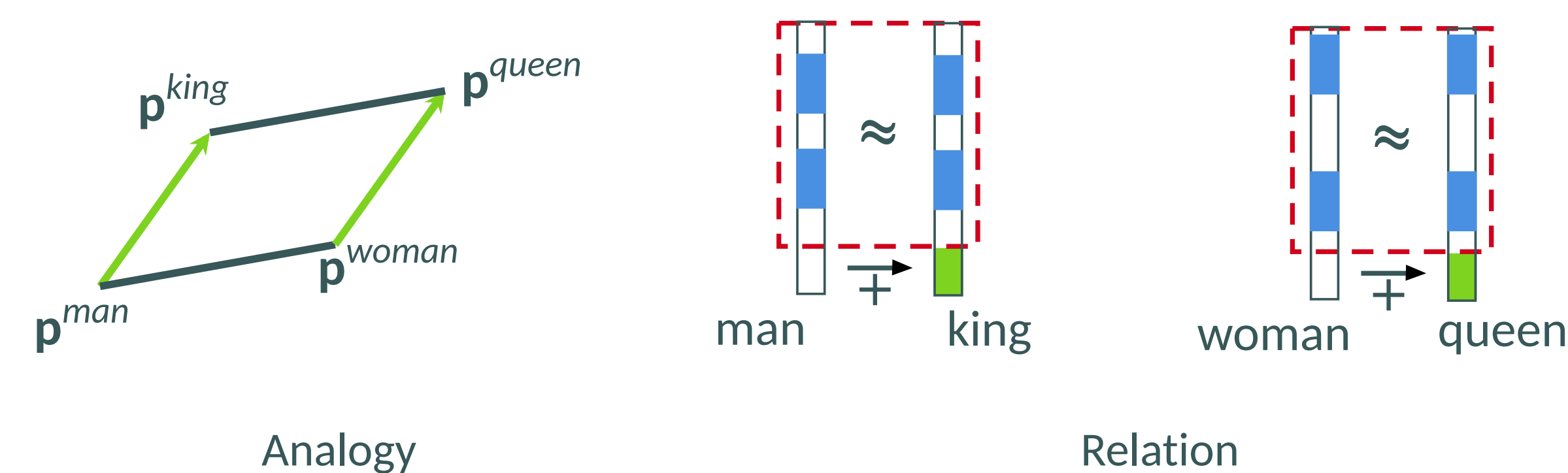
**Analogies:** word pairs that share a similar semantic difference, e.g. {man, king} and {woman, queen}. A linear combination of PMI vectors identifies analogies (subject to similar error terms, Allen and Hospedales (2019)):

$$p^{b*} - p^b \approx p^{a*} - p^a$$

All interactions *linear*

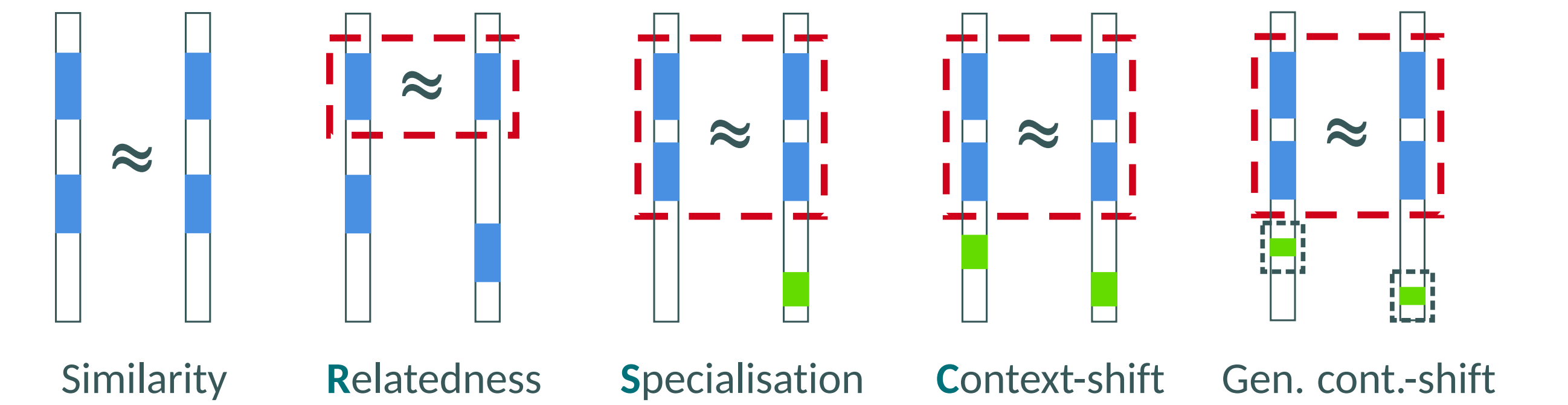
$\implies$  preserved in **word embeddings** under ( $\approx$ ) linear projection

## From Analogies to KG Relations



- Analogies contain common **binary word relations**, similar to KGs.
- For certain analogies ("specialisations"), the associated "vector offset" gives an **affine transformation** that represents the relation.
- Not all relations fit this semantic pattern, but we can now consider geometric aspects (**relation conditions**) of other relation types.

## Categorising Relations: Semantics $\rightarrow$ Relation Conditions



**Categorical completeness:** view PMI vectors as **sets** of word features and relation types as **set operations**.

- similarity  $\implies$  set equality
- relatedness  $\implies$  subset equality (relation-specific)
- (gen.) context-shift  $\implies$  set difference (relation-specific)

## Relation Conditions $\rightarrow$ Mappings between Embeddings

- R:**  $\mathcal{S}$ -relatedness requires  $e_s$  and  $e_o$  to share a subspace component  $\mathbb{V}_{\mathcal{S}}$ 
  - project onto  $\mathbb{V}_{\mathcal{S}}$  (multiply by matrix  $P_r \in \mathbb{R}^{d \times d}$ ) and compare.
  - Dot product:  $(P_r e_s)^T (P_r e_o) = e_s^T P_r^T P_r e_o = e_s^T M_r e_o$
  - Euclidean distance:  $\|P_r e_s - P_r e_o\|^2 = \|P_r e_s\|^2 - 2e_s^T M_r e_o + \|P_r e_o\|^2$
- S/C:** requires  $\mathcal{S}$ -relatedness and relation-specific component(s) ( $v_r^s, v_r^o$ ).
  - project onto a subspace corresponding to  $\mathcal{S}$ ,  $v_r^s$  and  $v_r^o$  (i.e. test  $\mathcal{S}$ -relatedness while preserving relation-specific components);
  - add relation-specific  $r = v_r^o - v_r^s \in \mathbb{R}^d$  to transformed embeddings.
  - Dot product:  $(P_r e_s + r)^T P_r e_o$
  - Euclidean distance:  $\|P_r e_s + r - P_r e_o\|^2$

## References

- Carl Allen and Timothy Hospedales. Analogies explained: Towards understanding word embeddings. In *International Conference on Machine Learning*, 2019.
- Carl Allen, Ivana Balažević, and Timothy Hospedales. What the Vec? Towards Probabilistically Grounded Embeddings. In *Advances in Neural Information Processing Systems*, 2019.
- Omer Levy and Yoav Goldberg. Neural word embedding as implicit matrix factorization. In *Advances in Neural Information Processing Systems*, 2014.
- 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*, 2013.