

Uncovering divergent linguistic information in word embeddings

VL Embeddings

Uni Heidelberg

SS 2019

Uncovering linguistic information in word embeddings

Artetxe et al (2018): Uncovering divergent linguistic information in word embeddings with lessons for intrinsic and extrinsic evaluation

- Word embeddings capture more information than we can directly observe
 - We can apply linear transformations to pretrained embeddings to adjust performance to different tasks along the axes of *similarity/relatedness* and *semantics/syntax*

Uncovering linguistic information in word embeddings

Artetxe et al (2018): Uncovering divergent linguistic information in word embeddings with lessons for intrinsic and extrinsic evaluation

- Word embeddings capture more information than we can directly observe
 - We can apply linear transformations to pretrained embeddings to adjust performance to different tasks along the axes of *similarity/relatedness* and *semantics/syntax*

	FastText	Dep-based embeddings	GloVe
good for	syntactic analogies	functional similarities	semantic analogies

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary
- The dot product $sim(i, j) = X_i \cdot X_j$ is a measure of the similarity between the i th and the j th word

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary
- The dot product $sim(i, j) = X_i \cdot X_j$ is a measure of the similarity between the i th and the j th word
- Define the similarity matrix $M(X) := XX^T$ so that $sim(i, j) = M(X)_{ij} \Rightarrow$ first order similarity

Linear transformation of embedding matrix

Artetxe et al (2018)

$$X = \begin{matrix} & \text{cat} & -0.19 & 0.45 & -0.40 \\ & \text{dog} & -0.28 & 0.43 & -0.39 \\ & \text{blue} & 0.02 & -0.40 & -0.39 \\ & \text{red} & 0.03 & -0.22 & -0.31 \\ & \text{happy} & -0.03 & -0.07 & -0.19 \end{matrix}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$X = \begin{matrix} & \text{cat} & -0.19 & 0.45 & -0.40 \\ & \text{dog} & -0.28 & 0.43 & -0.39 \\ & \text{blue} & 0.02 & -0.40 & -0.39 \\ & \text{red} & 0.03 & -0.22 & -0.31 \\ & \text{happy} & -0.03 & -0.07 & -0.19 \end{matrix}$$

$$X^T = \begin{matrix} & \text{cat} & \text{dog} & \text{blue} & \text{red} & \text{happy} \\ -0.19 & -0.28 & 0.02 & 0.03 & -0.03 \\ 0.45 & 0.43 & -0.40 & -0.22 & -0.07 \\ -0.40 & -0.39 & -0.39 & -0.31 & -0.19 \end{matrix}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$X = \begin{matrix} & \text{cat} & \text{dog} & \text{blue} & \text{red} & \text{happy} \\ \text{cat} & -0.19 & 0.45 & -0.40 & & \\ \text{dog} & -0.28 & 0.43 & -0.39 & & \\ \text{blue} & 0.02 & -0.40 & -0.39 & & \\ \text{red} & 0.03 & -0.22 & -0.31 & & \\ \text{happy} & -0.03 & -0.07 & -0.19 & & \end{matrix}$$

$$\text{sim}(i, j) = X_i \cdot X_j$$

$$X^T = \begin{matrix} & \text{cat} & \text{dog} & \text{blue} & \text{red} & \text{happy} \\ -0.19 & -0.28 & 0.02 & 0.03 & -0.03 & \\ 0.45 & 0.43 & -0.40 & -0.22 & -0.07 & \\ -0.40 & -0.39 & -0.39 & -0.31 & -0.19 & \end{matrix}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$X = \begin{matrix} & \text{cat} & -0.19 & 0.45 & -0.40 \\ & \text{dog} & -0.28 & 0.43 & -0.39 \\ & \text{blue} & 0.02 & -0.40 & -0.39 \\ & \text{red} & 0.03 & -0.22 & -0.31 \\ & \text{happy} & -0.03 & -0.07 & -0.19 \end{matrix}$$

$$\text{sim}(i, j) = X_i \cdot X_j$$

$$X^T = \begin{matrix} & \text{cat} & \text{dog} & \text{blue} & \text{red} & \text{happy} \\ -0.19 & -0.28 & 0.02 & 0.03 & -0.03 \\ 0.45 & 0.43 & -0.40 & -0.22 & -0.07 \\ -0.40 & -0.39 & -0.39 & -0.31 & -0.19 \end{matrix}$$

$$M(X) := XX^T$$

$$\text{sim}(i, j) = M(X)_{ij}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary
- The dot product $sim(i, j) = X_i \cdot X_j$ is a measure of the similarity between the i th and the j th word
- Define the similarity matrix $M(X) := XX^T$ so that $sim(i, j) = M(X)_{ij} \Rightarrow$ first order similarity

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary
- The dot product $sim(i, j) = X_i \cdot X_j$ is a measure of the similarity between the i th and the j th word
- Define the similarity matrix $M(X) := XX^T$ so that $sim(i, j) = M(X)_{ij} \Rightarrow$ first order similarity
- We can also define a second order similarity measure:
 - first order: How similar are w_i and w_j ?
 - second order: How similar are the contexts of w_i and w_j ?

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary
- The dot product $sim(i, j) = X_i \cdot X_j$ is a measure of the similarity between the i th and the j th word
- Define the similarity matrix $M(X) := XX^T$ so that $sim(i, j) = M(X)_{ij} \Rightarrow$ first order similarity
- We can also define a second order similarity measure:
 - first order: How similar are w_i and w_j ?
 - second order: How similar are the contexts of w_i and w_j ?
- We can even define a third, fourth or n th order similarity

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let X be the word embeddings matrix
- Let X_i be the embedding of the i th word in the vocabulary
- The dot product $sim(i, j) = X_i \cdot X_j$ is a measure of the similarity between the i th and the j th word
- Define the similarity matrix $M(X) := XX^T$ so that $sim(i, j) = M(X)_{ij} \Rightarrow$ first order similarity
- We can also define a second order similarity measure:
 - first order: How similar are w_i and w_j ?
 - second order: How similar are the contexts of w_i and w_j ?
- We can even define a third, fourth or n th order similarity

Idea: Some higher order similarities might be better at capturing specific aspects of language.

Linear transformation of embedding matrix

Artetxe et al (2018)

- Define the **first order** similarity matrix as

$$M(X) := XX^{\top} \quad \text{so that } sim(i, j) = M(X)_{ij}$$

- Define the **second order** similarity matrix as

$$M_2(X) := M(M(X)) \quad \text{so that } sim_2(i, j) = M_2(X)_{ij}$$

$$\text{where } M_2(X) = XX^{\top}XX^{\top}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

Definition: $M_2(X) := M(M(X))$

Linear transformation of embedding matrix

Artetxe et al (2018)

Definition: $M_2(X) := M(M(X))$

$$M_2(X) = M(M(X))$$

Linear transformation of embedding matrix

Artetxe et al (2018)

Definition: $M_2(X) := M(M(X))$

$$\begin{aligned}M_2(X) &= M(M(X)) \\ &= M(X)M(X)^\top\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

Definition: $M_2(X) := M(M(X))$

$$\begin{aligned}M_2(X) &= M(M(X)) \\ &= M(X)M(X)^\top \\ &= XX^\top(XX^\top)^\top\end{aligned}$$

$$(AB)^\top = B^\top A^\top$$

Linear transformation of embedding matrix

Artetxe et al (2018)

Definition: $M_2(X) := M(M(X))$

$$\begin{aligned}M_2(X) &= M(M(X)) \\ &= M(X)M(X)^\top \\ &= XX^\top(XX^\top)^\top \\ &= XX^\top X^{\top\top} X^\top\end{aligned}$$

$$(AB)^\top = B^\top A^\top$$

Linear transformation of embedding matrix

Artetxe et al (2018)

Definition: $M_2(X) := M(M(X))$

$$\begin{aligned}M_2(X) &= M(M(X)) \\ &= M(X)M(X)^\top \\ &= XX^\top(XX^\top)^\top \\ &= XX^\top X^{\top\top} X^\top \\ &= XX^\top XX^\top\end{aligned}$$

$$(AB)^\top = B^\top A^\top$$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Define the second order similarity matrix as

$$M_2(X) := XX^\top XX^\top \quad \text{so that } \text{sim}_2(i, j) = M_2(X)_{ij}$$

where $M_2(X) = M(M(X))$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Define the second order similarity matrix as
 $M_2(X) := XX^\top XX^\top$ so that $sim_2(i, j) = M_2(X)_{ij}$
where $M_2(X) = M(M(X))$
- Define the n -th order similarity matrix as
 $M_n(X) = (XX^\top)^n$ so that $sim_n(i, j) = M_n(X)_{ij}$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Define the second order similarity matrix as
 $M_2(X) := XX^\top XX^\top$ so that $sim_2(i, j) = M_2(X)_{ij}$
where $M_2(X) = M(M(X))$
- Define the n -th order similarity matrix as
 $M_n(X) = (XX^\top)^n$ so that $sim_n(i, j) = M_n(X)_{ij}$

Instead of changing the similarity measure, we can also change the word embeddings themselves through a linear transformation so they directly capture this second or n -th order similarity

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let $X^T X = Q \Lambda Q^T$ be the eigendecomposition of $X^T X$
 - Λ is a positive diagonal matrix whose entries are the eigenvalues of $X^T X$ and Q is an orthogonal matrix with their respective eigenvectors as columns

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let $X^T X = Q \Lambda Q^T$ be the eigendecomposition of $X^T X$
 - Λ is a positive diagonal matrix whose entries are the eigenvalues of $X^T X$ and Q is an orthogonal matrix with their respective eigenvectors as columns

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X \\ m \times n \end{matrix} = \begin{matrix} U \\ m \times m \end{matrix} \cdot \begin{matrix} \Sigma \\ m \times n \end{matrix} \cdot \begin{matrix} Q^T \\ n \times n \end{matrix}$$

SVD

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X & = & U & \cdot & \Sigma & \cdot & Q^T \\ m \times n & & m \times m & & m \times n & & n \times n \end{matrix} \quad \text{SVD}$$

$$X^T X = (U \cdot \Sigma \cdot Q^T)^T \cdot (U \cdot \Sigma \cdot Q^T)$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X \\ m \times n \end{matrix} = \begin{matrix} U & \cdot & \Sigma & \cdot & Q^T \\ m \times m & & m \times n & & n \times n \end{matrix} \quad \text{SVD}$$

$$\begin{aligned} X^T X &= (U \cdot \Sigma \cdot Q^T)^T \cdot (U \cdot \Sigma \cdot Q^T) \\ &= (Q^T)^T \cdot \Sigma^T \cdot U^T \cdot U \cdot \Sigma \cdot Q^T \end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X \\ m \times n \end{matrix} = \begin{matrix} U & \cdot & \Sigma & \cdot & Q^T \\ m \times m & & m \times n & & n \times n \end{matrix} \quad \text{SVD}$$

$$\begin{aligned} X^T X &= (U \cdot \Sigma \cdot Q^T)^T \cdot (U \cdot \Sigma \cdot Q^T) \\ &= (Q^T)^T \cdot \Sigma^T \cdot U^T \cdot U \cdot \Sigma \cdot Q^T \quad U \text{ orthonormal} \end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X \\ m \times n \end{matrix} = \begin{matrix} U & \cdot & \Sigma & \cdot & Q^T \\ m \times m & & m \times n & & n \times n \end{matrix} \quad \text{SVD}$$

$$\begin{aligned} X^T X &= (U \cdot \Sigma \cdot Q^T)^T \cdot (U \cdot \Sigma \cdot Q^T) \\ &= (Q^T)^T \cdot \Sigma^T \cdot U^T \cdot U \cdot \Sigma \cdot Q^T \quad U \text{ orthonormal} \\ &= Q \cdot \Sigma^T \cdot \Sigma \cdot Q^T \end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X \\ m \times n \end{matrix} = \begin{matrix} U & \cdot & \Sigma & \cdot & Q^T \\ m \times m & & m \times n & & n \times n \end{matrix} \quad \text{SVD}$$

$$\begin{aligned} X^T X &= (U \cdot \Sigma \cdot Q^T)^T \cdot (U \cdot \Sigma \cdot Q^T) \\ &= (Q^T)^T \cdot \Sigma^T \cdot U^T \cdot U \cdot \Sigma \cdot Q^T \quad U \text{ orthonormal} \\ &= Q \cdot \Sigma^T \cdot \Sigma \cdot Q^T \quad \Sigma \text{ Diagonalmatrix mit } \sqrt{EW} \end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{matrix} X \\ m \times n \end{matrix} = \begin{matrix} U & \cdot & \Sigma & \cdot & Q^T \\ m \times m & & m \times n & & n \times n \end{matrix} \quad \text{SVD}$$

$$\begin{aligned} X^T X &= (U \cdot \Sigma \cdot Q^T)^T \cdot (U \cdot \Sigma \cdot Q^T) \\ &= (Q^T)^T \cdot \Sigma^T \cdot U^T \cdot U \cdot \Sigma \cdot Q^T \quad U \text{ orthonormal} \\ &= Q \cdot \Sigma^T \cdot \Sigma \cdot Q^T \quad \Sigma \text{ Diagonalmatrix mit } \sqrt{EW} \\ &= Q \cdot \Lambda \cdot Q^T \end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let $X^T X = Q \Lambda Q^T$ be the eigendecomposition of $X^T X$
 - Λ is a positive diagonal matrix whose entries are the eigenvalues of $X^T X$ and Q is an orthogonal matrix with their respective eigenvectors as columns

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let $X^T X = Q \Lambda Q^T$ be the eigendecomposition of $X^T X$
 - Λ is a positive diagonal matrix whose entries are the eigenvalues of $X^T X$ and Q is an orthogonal matrix with their respective eigenvectors as columns
- Define the **linear transformation matrix** $W := Q\sqrt{\Lambda}$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let $X^T X = Q \Lambda Q^T$ be the eigendecomposition of $X^T X$
 - Λ is a positive diagonal matrix whose entries are the eigenvalues of $X^T X$ and Q is an orthogonal matrix with their respective eigenvectors as columns
- Define the **linear transformation matrix** $W := Q\sqrt{\Lambda}$
- Apply W to the original embeddings $X \Rightarrow X' = XW$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Let $X^T X = Q \Lambda Q^T$ be the eigendecomposition of $X^T X$
 - Λ is a positive diagonal matrix whose entries are the eigenvalues of $X^T X$ and Q is an orthogonal matrix with their respective eigenvectors as columns
- Define the **linear transformation matrix** $W := Q \sqrt{\Lambda}$
- Apply W to the original embeddings $X \Rightarrow X' = XW$
- $M(X') = M_2(X) \Rightarrow$ transformed embeddings X' capture the **second order** similarity as defined for the original embeddings

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^{\top}$$

first order similarity of X'

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(X') &= X' \cdot X'^{\top} && \text{first order similarity of } X' \\ &= X \cdot W \cdot (XW)^{\top}\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^{\top}$$

first order similarity of X'

$$= X \cdot W \cdot (XW)^{\top}$$

$$(AB)^{\top} = B^{\top} A^{\top}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(X') &= X' \cdot X'^{\top} && \text{first order similarity of } X' \\ &= X \cdot W \cdot (XW)^{\top} && (AB)^{\top} = B^{\top} A^{\top} \\ &= X \cdot W \cdot W^{\top} \cdot X^{\top}\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^{\top} \quad \text{first order similarity of } X'$$

$$= X \cdot W \cdot (XW)^{\top} \quad (AB)^{\top} = B^{\top} A^{\top}$$

$$= X \cdot W \cdot W^{\top} \cdot X^{\top} \quad W := Q\sqrt{\Lambda}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^T \quad \text{first order similarity of } X'$$

$$= X \cdot W \cdot (XW)^T \quad (AB)^T = B^T A^T$$

$$= X \cdot W \cdot W^T \cdot X^T \quad W := Q\sqrt{\Lambda}$$

$$= X \cdot Q\sqrt{\Lambda} \cdot (Q\sqrt{\Lambda})^T \cdot X^T$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^T \quad \text{first order similarity of } X'$$

$$= X \cdot W \cdot (XW)^T \quad (AB)^T = B^T A^T$$

$$= X \cdot W \cdot W^T \cdot X^T \quad W := Q\sqrt{\Lambda}$$

$$= X \cdot Q\sqrt{\Lambda} \cdot (Q\sqrt{\Lambda})^T \cdot X^T \quad (AB)^T = B^T A^T$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^T \quad \text{first order similarity of } X'$$

$$= X \cdot W \cdot (XW)^T \quad (AB)^T = B^T A^T$$

$$= X \cdot W \cdot W^T \cdot X^T \quad W := Q\sqrt{\Lambda}$$

$$= X \cdot Q\sqrt{\Lambda} \cdot (Q\sqrt{\Lambda})^T \cdot X^T \quad (AB)^T = B^T A^T$$

$$= X \cdot Q\sqrt{\Lambda} \cdot \sqrt{\Lambda}^T \cdot Q^T X^T$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(X') &= X' \cdot X'^T && \text{first order similarity of } X' \\&= X \cdot W \cdot (XW)^T && (AB)^T = B^T A^T \\&= X \cdot W \cdot W^T \cdot X^T && W := Q\sqrt{\Lambda} \\&= X \cdot Q\sqrt{\Lambda} \cdot (Q\sqrt{\Lambda})^T \cdot X^T && (AB)^T = B^T A^T \\&= X \cdot Q\sqrt{\Lambda} \cdot \sqrt{\Lambda}^T \cdot Q^T X^T \\&= X \cdot Q \cdot \Lambda \cdot Q^T \cdot X^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(X') &= X' \cdot X'^T && \text{first order similarity of } X' \\ &= X \cdot W \cdot (XW)^T && (AB)^T = B^T A^T \\ &= X \cdot W \cdot W^T \cdot X^T && W := Q\sqrt{\Lambda} \\ &= X \cdot Q\sqrt{\Lambda} \cdot (Q\sqrt{\Lambda})^T \cdot X^T && (AB)^T = B^T A^T \\ &= X \cdot Q\sqrt{\Lambda} \cdot \sqrt{\Lambda}^T \cdot Q^T X^T \\ &= X \cdot Q \cdot \Lambda \cdot Q^T \cdot X^T \\ &= X \cdot X^T \cdot X \cdot X^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(X') = X' \cdot X'^T \quad \text{first order similarity of } X'$$

$$= X \cdot W \cdot (XW)^T \quad (AB)^T = B^T A^T$$

$$= X \cdot W \cdot W^T \cdot X^T \quad W := Q\sqrt{\Lambda}$$

$$= X \cdot Q\sqrt{\Lambda} \cdot (Q\sqrt{\Lambda})^T \cdot X^T \quad (AB)^T = B^T A^T$$

$$= X \cdot Q\sqrt{\Lambda} \cdot \sqrt{\Lambda}^T \cdot Q^T X^T$$

$$= X \cdot Q \cdot \Lambda \cdot Q^T \cdot X^T$$

$$= X \cdot X^T \cdot X \cdot X^T = M_2(X) \quad \text{second order similarity of } X$$

Linear transformation of embedding matrix

Artetxe et al (2018)

More generally

- Define $W_\alpha := Q\Lambda^\alpha$

where α is a parameter of the transformation that adjusts to the desired similarity order:

<i>first order similarity</i>	$\alpha = 0$	\Rightarrow	$M(XW_0) = M(X)$
<i>second order similarity</i>	$\alpha = 0.5$	\Rightarrow	$M(XW_{0.5}) = M_2(X)$
<i>n-th order similarity</i>	$\alpha = (n - 1)/2$	\Rightarrow	$M(XW_\alpha) = M_n(X)$

Linear transformation of embedding matrix

Artetxe et al (2018)

More generally

- Define $W_\alpha := Q\Lambda^\alpha$

where α is a parameter of the transformation that adjusts to the desired similarity order:

<i>first order similarity</i>	$\alpha = 0$	\Rightarrow	$M(XW_0) = M(X)$
<i>second order similarity</i>	$\alpha = 0.5$	\Rightarrow	$M(XW_{0.5}) = M_2(X)$
<i>n-th order similarity</i>	$\alpha = (n - 1)/2$	\Rightarrow	$M(XW_\alpha) = M_n(X)$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(XW_0) = M(X \cdot Q \cdot \Lambda^0) \quad \Lambda^0 \text{ Einheitsmatrix}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_0) &= M(X \cdot Q \cdot \Lambda^0) \quad \Lambda^0 \text{ Einheitsmatrix} \\ &= M(X \cdot Q)\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_0) &= M(X \cdot Q \cdot \Lambda^0) \quad \Lambda^0 \text{ Einheitsmatrix} \\ &= M(X \cdot Q) \\ &= X \cdot Q(XQ)^\top\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_0) &= M(X \cdot Q \cdot \Lambda^0) \quad \Lambda^0 \text{ Einheitsmatrix} \\ &= M(X \cdot Q) \\ &= X \cdot Q(XQ)^\top \\ &= X \cdot Q \cdot Q^\top \cdot X^\top\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_0) &= M(X \cdot Q \cdot \Lambda^0) \quad \Lambda^0 \text{ Einheitsmatrix} \\ &= M(X \cdot Q) \\ &= X \cdot Q(XQ)^\top \\ &= X \cdot Q \cdot Q^\top \cdot X^\top \\ &= X \cdot X^\top\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_0) &= M(X \cdot Q \cdot \Lambda^0) \quad \Lambda^0 \text{ Einheitsmatrix} \\ &= M(X \cdot Q) \\ &= X \cdot Q(XQ)^\top \\ &= X \cdot Q \cdot Q^\top \cdot X^\top \\ &= X \cdot X^\top \\ &= M(X)\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

More generally

- Define $W_\alpha := Q\Lambda^\alpha$

where α is a parameter of the transformation that adjusts to the desired similarity order:

first order similarity $\alpha = 0$ \Rightarrow $M(XW_0) = M(X)$

second order similarity $\alpha = 0.5$ \Rightarrow $M(XW_{0.5}) = M_2(X)$

n-th order similarity $\alpha = (n - 1)/2$ \Rightarrow $M(XW_\alpha) = M_n(X)$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$M(XW_{0.5}) = M(X \cdot Q\sqrt{\Lambda})$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_{0.5}) &= M(X \cdot Q\sqrt{\Lambda}) \\ &= X \cdot Q\sqrt{\Lambda} (X \cdot Q\sqrt{\Lambda})^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_{0.5}) &= M(X \cdot Q\sqrt{\Lambda}) \\ &= X \cdot Q\sqrt{\Lambda} (X \cdot Q\sqrt{\Lambda})^T \\ &= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda}^T Q^T \cdot X^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_{0.5}) &= M(X \cdot Q\sqrt{\Lambda}) \\ &= X \cdot Q\sqrt{\Lambda} (X \cdot Q\sqrt{\Lambda})^T \\ &= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda}^T Q^T \cdot X^T \\ &= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda} \cdot Q^T \cdot X^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_{0.5}) &= M(X \cdot Q\sqrt{\Lambda}) \\&= X \cdot Q\sqrt{\Lambda} (X \cdot Q\sqrt{\Lambda})^T \\&= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda}^T Q^T \cdot X^T \\&= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda} \cdot Q^T \cdot X^T \\&= X \cdot Q \cdot \Lambda \cdot Q^T \cdot X^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_{0.5}) &= M(X \cdot Q\sqrt{\Lambda}) \\&= X \cdot Q\sqrt{\Lambda} (X \cdot Q\sqrt{\Lambda})^T \\&= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda}^T Q^T \cdot X^T \\&= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda} \cdot Q^T \cdot X^T \\&= X \cdot Q \cdot \Lambda \cdot Q^T \cdot X^T \\&= X \cdot X^T \cdot X \cdot X^T\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

$$\begin{aligned}M(XW_{0.5}) &= M(X \cdot Q\sqrt{\Lambda}) \\&= X \cdot Q\sqrt{\Lambda} (X \cdot Q\sqrt{\Lambda})^\top \\&= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda}^\top Q^\top \cdot X^\top \\&= X \cdot Q\sqrt{\Lambda} \sqrt{\Lambda} \cdot Q^\top \cdot X^\top \\&= X \cdot Q \cdot \Lambda \cdot Q^\top \cdot X^\top \\&= X \cdot X^\top \cdot X \cdot X^\top \\&= M_2(X)\end{aligned}$$

Linear transformation of embedding matrix

Artetxe et al (2018)

More generally

- Define $W_\alpha := Q\Lambda^\alpha$

where α is a parameter of the transformation that adjusts to the desired similarity order:

<i>first order similarity</i>	$\alpha = 0$	\Rightarrow	$M(XW_0) = M(X)$
<i>second order similarity</i>	$\alpha = 0.5$	\Rightarrow	$M(XW_{0.5}) = M_2(X)$
<i>n-th order similarity</i>	$\alpha = (n - 1)/2$	\Rightarrow	$M(XW_\alpha) = M_n(X)$

Linear transformation of embedding matrix

Artetxe et al (2018)

- Assuming that the embeddings X capture some second order similarity, it is possible to transform them so that they capture the corresponding first order similarity
 - One can easily generalise this to higher order similarities by using smaller values of α
- ⇒ Parameter α can be used to either increase or decrease the similarity order that we want our embeddings to capture
- ⇒ α can be continuous

Linear transformation of different embeddings

Artetxe et al (2018)

		Word analogy		Word similarity	
		Semantic	Syntactic	Similarity (SimLex-999)	Relatedness (MEN)
word2vec	Original	76.49	74.87	44.21	76.96
	Best				
glove	Original	83.17	76.19	40.70	80.06
	Best				
fasttext	Original	89.76	82.44	50.48	83.55
	Best				

Table 1: Results in intrinsic evaluation for the original embeddings and the best post-processed model with the corresponding value of α . The evaluation measure is accuracy for word analogy and Spearman correlation for word similarity.

Linear transformation of different embeddings

Artetxe et al (2018)

		Word analogy		Word similarity	
		Semantic	Syntactic	Similarity (SimLex-999)	Relatedness (MEN)
word2vec	Original	76.49	74.87	44.21	76.96
	Best	81.00 $\alpha = -0.65$	74.96 $\alpha = 0.10$	47.81 $\alpha = -0.70$	78.09 $\alpha = -0.30$
glove	Original	83.17	76.19	40.70	80.06
	Best	86.73 $\alpha = -0.85$	76.51 $\alpha = -0.10$	51.54 $\alpha = -0.85$	84.00 $\alpha = -0.45$
fasttext	Original	89.76	82.44	50.48	83.55
	Best	90.85 $\alpha = -0.45$	84.45 $\alpha = 0.25$	51.55 $\alpha = -0.25$	84.06 $\alpha = -0.15$

Table 1: Results in intrinsic evaluation for the original embeddings and the best post-processed model with the corresponding value of α . The evaluation measure is accuracy for word analogy and Spearman correlation for word similarity.

Lessons learned for intrinsic and extrinsic evaluations

Artetxe et al (2018)

- Standard intrinsic evaluation is static and incomplete
 - ⇒ Intrinsic evaluation not a good predictor for performance in downstream applications
 - ⇒ Systems that use embeddings as features can learn task-specific optimal balance between the two axes

References

- Mikolov, Yih and Zweig: (2013): Linguistic regularities in continuous space word representations. NAACL 2013.
- Faruqui, Tsvetkov, Rastogi and Dyer (2016): Problems With Evaluation of Word Embeddings Using Word Similarity Tasks. The 1st Workshop on Evaluating Vector Space Representations for NLP, Berlin, Germany.
- Artetxe, Labaka, Lopez-Gazpio and Agirre (2018): Uncovering Divergent Linguistic Information in Word Embeddings with Lessons for Intrinsic and Extrinsic Evaluation. CoNLL 2018. Brussels, Belgium.
- Rubenstein and Goodenough (1965): Contextual correlates of synonymy. Communications of the ACM 8(10):627–633.
- Harris, Z. (1954). Distributional structure. Word, 10(23): 146-162.
- Multimodal Distributional Semantics E. Bruni, N. K. Tran and M. Baroni. Journal of Artificial Intelligence Research 49: 1-47.
- Collobert, Weston Bottou, Karlen, Kavukcuoglu and Kuksa (2011): Natural Language Processing (almost) from Scratch. Journal of Machine Learning Research 12 (2011) 2461-2505.
- Lu, Wang, Bansal, Gimpel and Livescu (2015): Deep multilingual correlation for improved word embeddings. NAACL 2015.
- Rastogi, Van Durme and Arora (2015): Multiview LSA: Representation learning via generalized CCA. NAACL 2015.
- Chiu, Korhonen and Pyysalo (2016): Intrinsic Evaluation of Word Vectors Fails to Predict Extrinsic Performance. ACL 2016.
- Data and Code
 - Code for Artetxe et al. (2018): <https://github.com/artetxem/uncovec>
 - The MEN dataset <https://staff.fnwi.uva.nl/e.bruni/MEN>
 - Datasets for word vector evaluation <https://github.com/vector-ai/word-benchmarks>