



Correlations between sets of Words

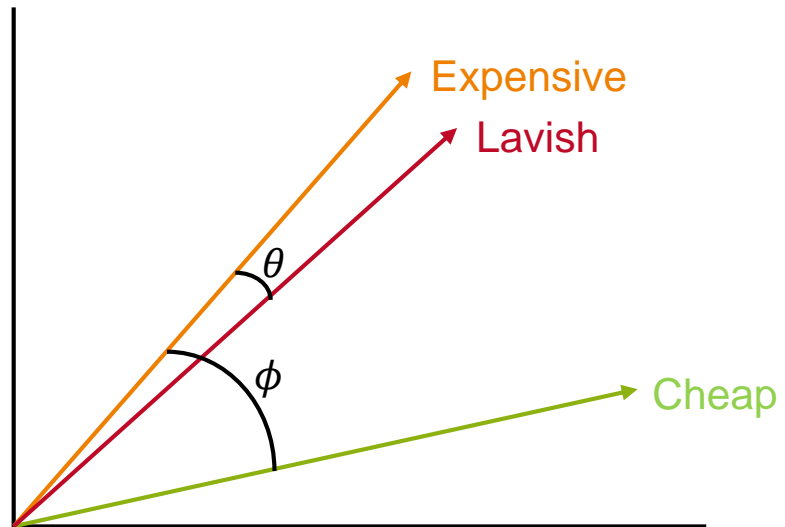
Carel van Niekerk

03.07.2020

1. Problem description
2. Two different views on word distributions
3. Standard Similarity measures
4. Similarity measures for sets
5. Reproducing Kernel Hilbert Spaces
6. Correlation between sets of words
7. Applications in Dialogue

Problem Description

- Given word embeddings for some words we can find semantic similarity between these words using cosine similarity.



- Now we have two sets of words:

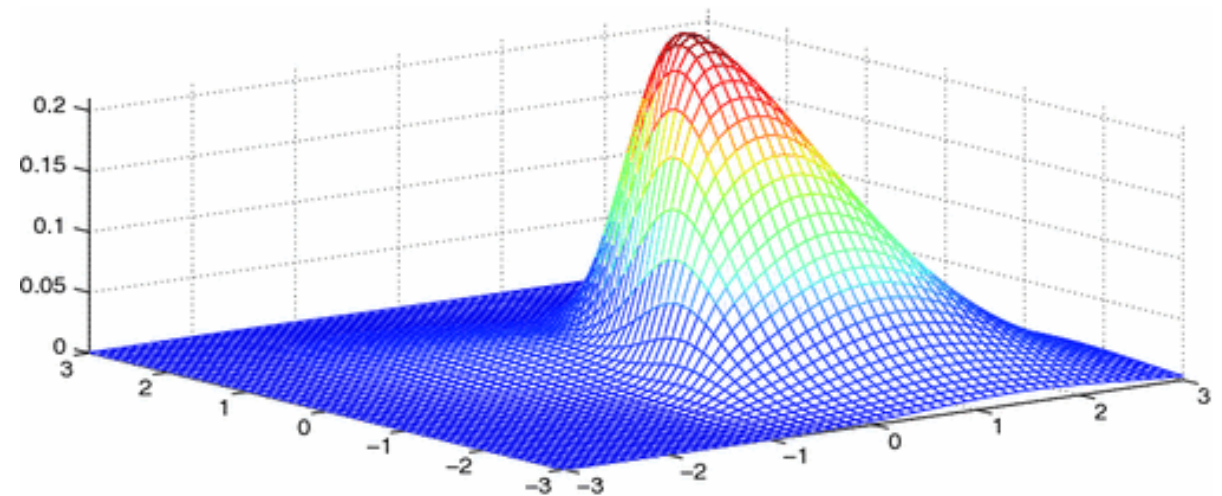
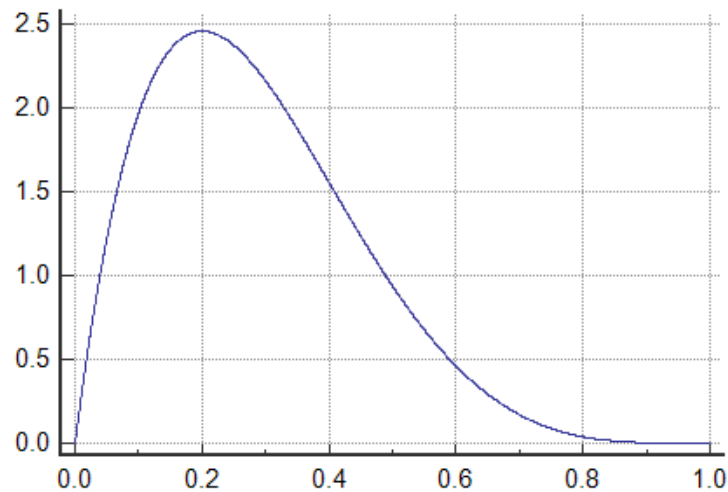
Not to pricey

Costing an arm and a leg

- Embedding models will now give us two sets of embeddings
- How do we find the similarity between these two sets?

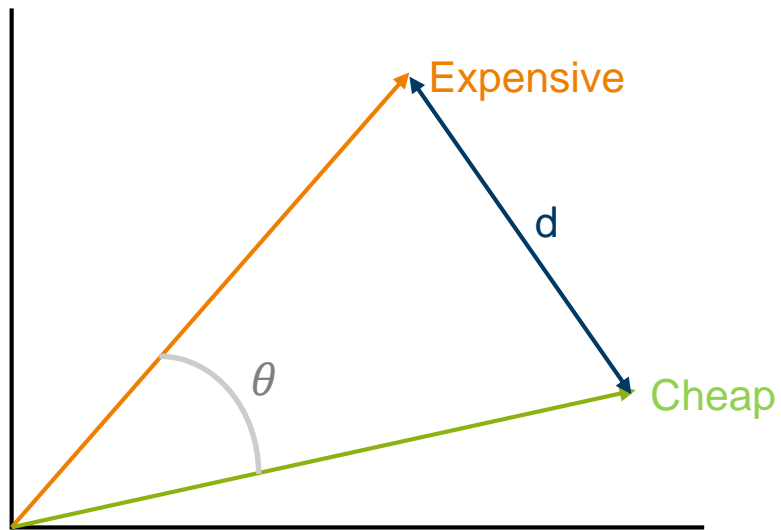
Two different views on word distributions

1. Each word embedding comes from a **d -dimensional distribution**.
2. Each word has its **own semantic distribution** (1-dimensional). The word embedding of that word is a set of observations from this distribution.



$$W_i \sim F_D(\cdot)$$

- Each word embedding is an observation from the d -dimensional distribution over semantic meanings.
- Similarity Measures:
 - Cosine Similarity
 - Euclidean Distance
- Cannot use correlation since we do not know the true distribution and we only have 1 observation of each word so it cannot be estimated.



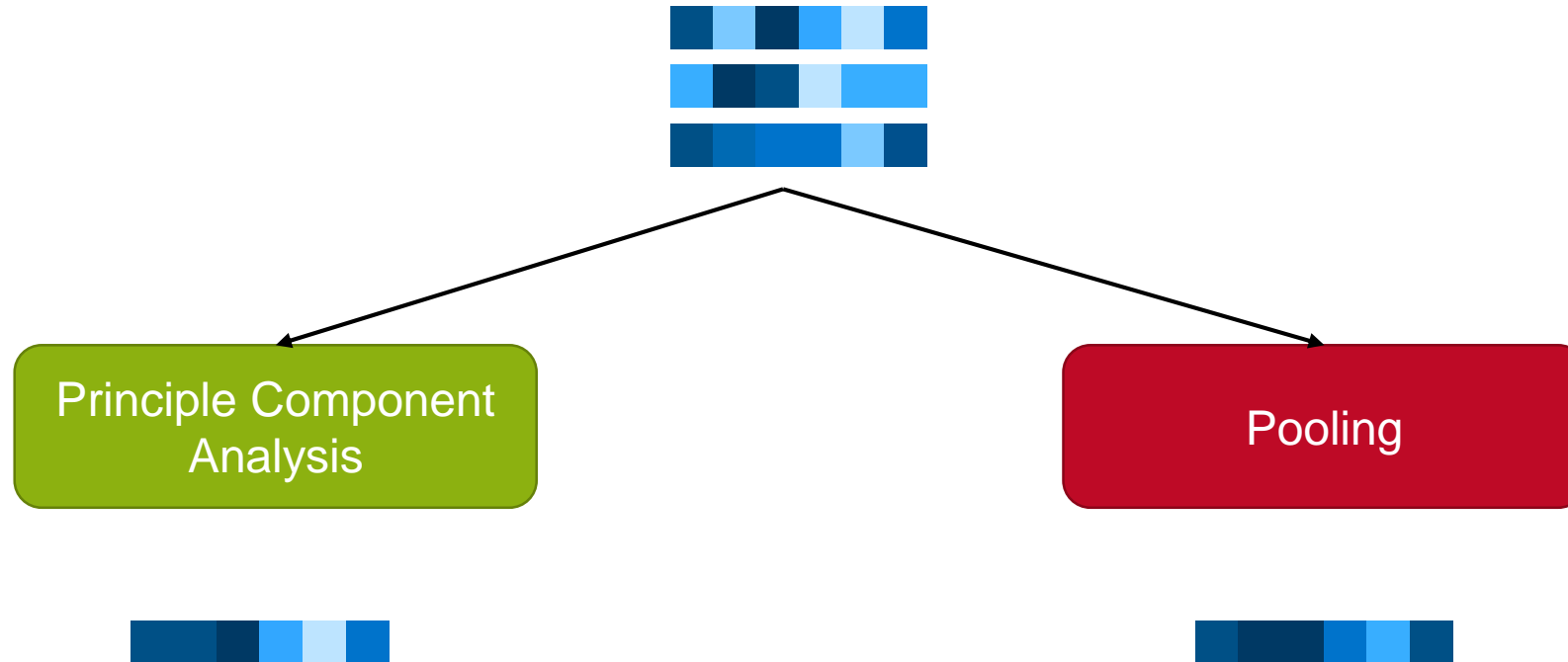
- Cosine Similarity

$$\text{Sim}(\text{Expensive}, \text{Cheap}) = \cos(\theta)$$

- Euclidean Distance (L2 Norm)

$$\text{Sim}(\text{Expensive}, \text{Cheap}) = d$$

Similarity of sets of words



Similarity of sets of words

Approach	STS 12	STS 13	STS 14	STS 15	STS 16
PCA	58.6	67.3	70.5	73.5	71.7
Mean Pool	58.3	57.9	64.9	67.6	64.3
Max Pool	57.7	53.5	67.2	69.5	68.5

Correlation between human annotated similarity and predicted similarity

Distribution of a word

$$W_i \sim F_i(\cdot)$$

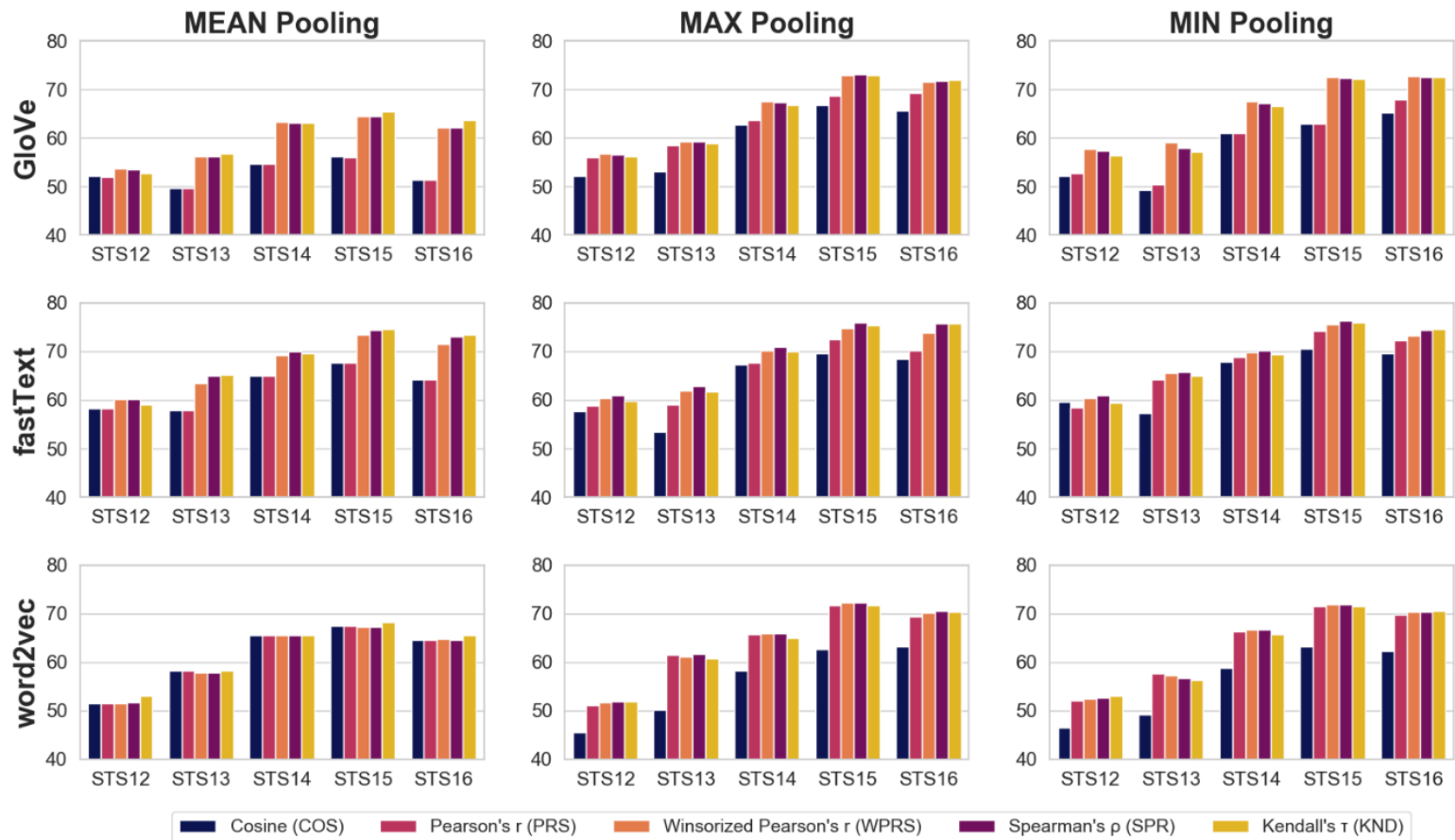
$$\mathbf{w}_i \in \mathbb{R}^D$$

- Each dimension in the word vector is an observation from the word distribution over semantic meanings.
- Similarity Measures:
 - Correlation Coefficient

■ Correlation

$$\begin{aligned} & \text{Sim}(\text{Expensive}, \text{Cheap}) \\ &= \text{Corr}(W_{\text{Expensive}}, W_{\text{Cheap}}) \\ &\approx \frac{\hat{\sigma}_{\text{Expensive, Cheap}}}{\hat{\sigma}_{\text{Expensive}} \hat{\sigma}_{\text{Cheap}}} \quad (\text{Estimates based on the observed word vectors}) \end{aligned}$$

If the average of the word vectors is 0, which is the case for models such as BERT, then we have that the correlation is equal to the cosine similarity.



- Let \mathbf{X} be an Arbitrary set and \mathbf{H} a Reproducing Kernel Hilbert Space with kernel \mathbf{K} .
- \mathbf{H} is a space of function over \mathbf{X} .
- Then for all x in \mathbf{X} there exists a unique function K_x in \mathbf{H} with the reproducing property $f(x) = \langle f, K_x \rangle_H$ for all f in \mathbf{H} .
- $\langle \cdot, \cdot \rangle_H$ is the inner product under \mathbf{H} .
- It has been shown that RKHS can be used to find the distance between two sets using only their Kernel functions. This is due to the reproducing property.

- Consider two sets of word embeddings $S_1 = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\}$, $S_2 = \{\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_M\}$
- We construct two RKHS \mathcal{F} and \mathcal{G} over these two sets.
- Then for kernel functions K and L we can calculate the Hilbert-Schmidt independence criterion (HSIC) which indicates the level of dependence between the two sets.
- For $\mathbf{K} = K(\mathbf{x}_i, \mathbf{x}_j)$ and $\mathbf{L} = L(\mathbf{y}_i, \mathbf{y}_j)$:

$$HSIC(S_1, S_2, K, L) = HSIC(\mathbf{K}, \mathbf{L})$$

- A generalisation of the Pearson's correlation coefficient for these sets is the Centred Kernel Alignment (CKA). The CKA is defined as:

$$CKA(\mathbf{K}, \mathbf{L}) = \frac{HSIC(\mathbf{K}, \mathbf{L})}{\sqrt{HSIC(\mathbf{K}, \mathbf{K})HSIC(\mathbf{L}, \mathbf{L})}}$$

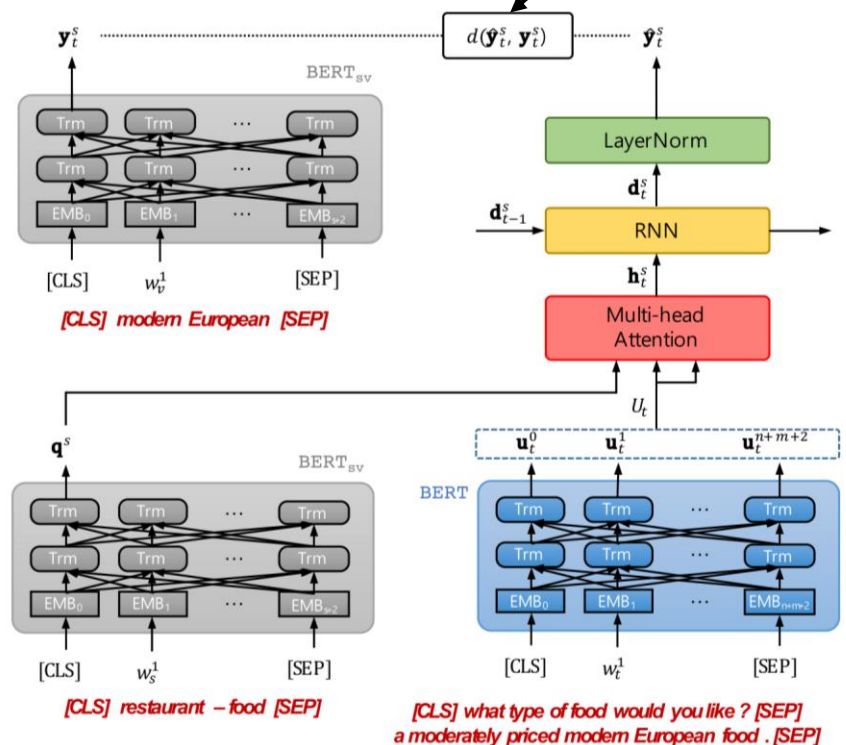
- HSIC is an operator comparable to Covariance and hence CKA is comparable to the Pearson's correlation coefficient.

Approach	STS 12	STS 13	STS 14	STS 15	STS 16
PCA	58.6	67.3	70.5	73.5	71.7
MeanPool	58.3	57.9	64.9	67.6	64.3
MaxPool	57.7	53.5	67.2	69.5	68.5
Linear CKA	59.8	62.1	69.5	74.6	70.3
Gaussian CKA	60.5	63.8	71.6	76.3	73.7
dCor CKA	61.0	63.2	71.5	75.6	72.4

Correlation between human annotated similarity and predicted similarity

- Semantic Similarity is a very useful tool in striving towards ontology independent Dialogue State Tracking.
 - Additional useful information and descriptions come in the form of sets of words.
- Semantic Similarity can be a useful tool in evaluating responses.
 - Responses are also sets of words
- This correlation between sets of words measure can be a useful tool in dialogue modelling.

Similarity between the predicted value and the list of possible values



- Currently a pooled approach with cosine similarity is used.
- This could potentially be replaced by the improved CKA correlation.

- Statistical correlation between sets of words promises to be a better estimate of semantic similarity than cosine similarity based on pooled / PCA vectors.
- Good semantic similarity measures combined with good embedding models could possibly improve ontology independent dialogue models.



Questions

- [A kernel method for the two-sample problem. Gretton et.al. 2007](#)
- [Correlations between Word Vector Sets. Zhelezniak et.al. 2019](#)
- [How Contextual are Contextualized Word Representations? Comparing the Geometry of BERT, ELMo, and GPT-2 Embeddings. Ethayarajh 2019](#)