




Spherical Text Embedding

Yu Meng¹, Jiaxin Huang¹, Guangyuan Wang¹, Chao Zhang²,
Honglei Zhuang³, Lance Kaplan⁴ and Jiawei Han¹

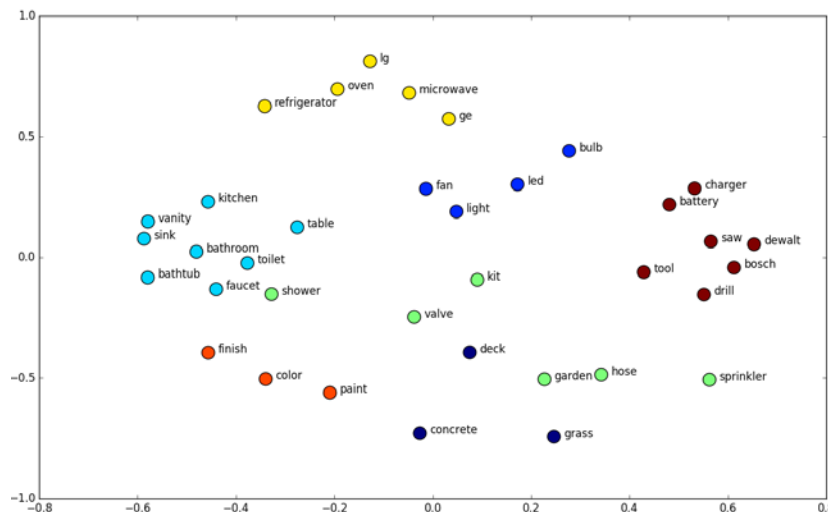
University of Illinois at Urbana-Champaign¹, Georgia Institute of Technology²,
Google Research³, U.S. Army Research Laboratory⁴

Outline

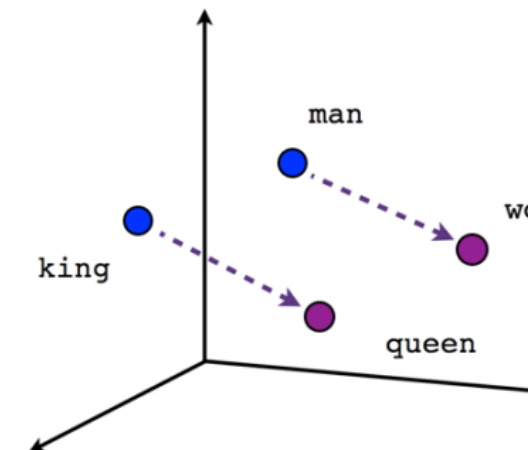
- ☐ Preliminaries 
- ☐ Motivations & Introduction
- ☐ Model: Spherical Text Embedding
- ☐ Method: Optimization on the Sphere
- ☐ Experiments
- ☐ Conclusions

Preliminaries: Text Embedding

- A milestone in NLP and ML:
 - Unsupervised learning of text representations—No supervision needed
 - Embed one-hot vectors into lower-dimensional space—Address “curse of dimensionality”
 - Word embedding captures useful properties of word semantics
 - Word similarity: Words with similar meanings are embedded closer
 - Word analogy: Linear relationships between words (e.g. king - queen = man - woman)



Word Similarity

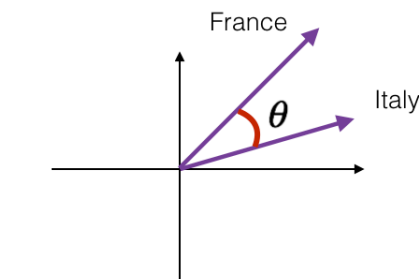


Word Analogy

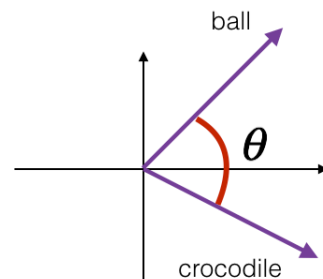
Preliminaries: Text Embedding

□ How to use text embeddings? Mostly directional similarity (i.e., cosine similarity)

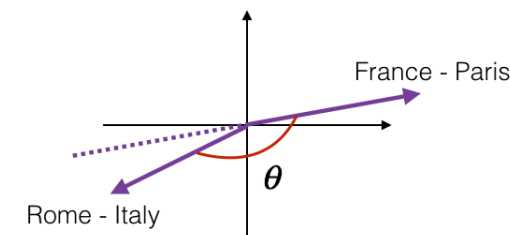
□ Word similarity is derived using cosine similarity



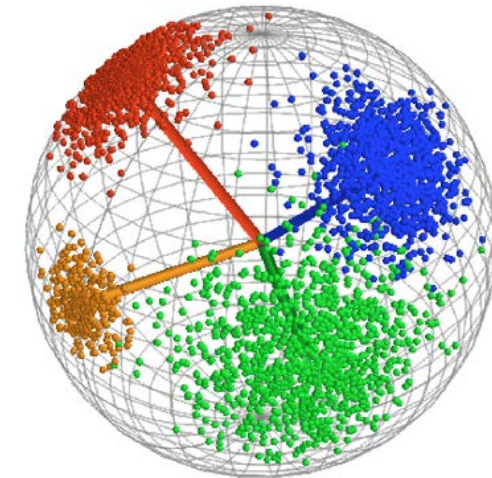
France and Italy are quite similar
 θ is close to 0°
 $\cos(\theta) \approx 1$



ball and crocodile are not similar
 θ is close to 90°
 $\cos(\theta) \approx 0$




the two vectors are similar but opposite
the first one encodes (city - country)
while the second one encodes (country - city)
 θ is close to 180°
 $\cos(\theta) \approx -1$



□ Word clustering (e.g. TaxoGen) is performed on a sphere

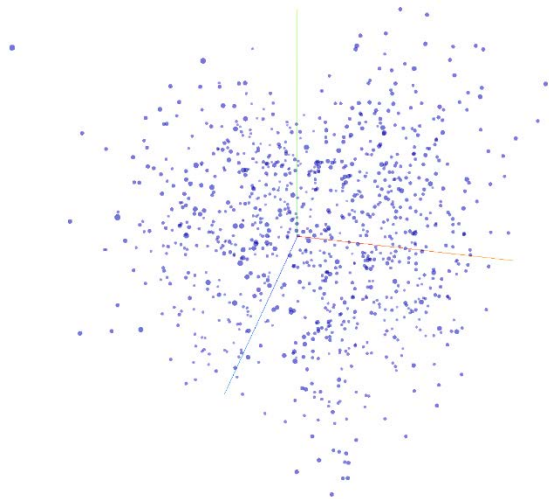
□ Better document clustering performances when embeddings are normalized and spherical clustering algorithms are used

Outline

- ❑ Preliminaries
- ❑ Motivations & Introduction 
- ❑ Model: Spherical Text Embedding
- ❑ Method: Optimization on the Sphere
- ❑ Experiments
- ❑ Conclusions

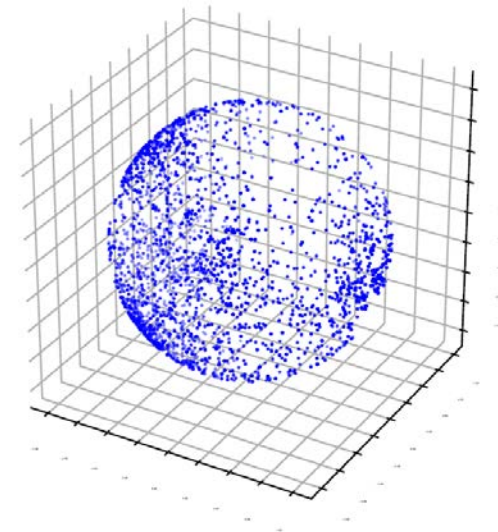
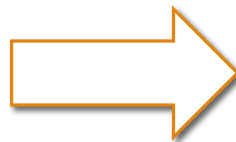
Motivations

- ❑ Issues with previous word embedding frameworks:
 - ❑ Although directional similarity has shown effective for various applications, previous embeddings (e.g. Word2Vec, GloVe, fastText) are trained in the Euclidean space
 - ❑ A gap between training space and usage space: Trained in Euclidean space but used on sphere



Embedding Training in Euclidean Space

Post-processing
(Normalization)



Embedding Usage on the Sphere
(Similarity, Clustering, etc.)

Motivations

- What is the consequence of the inconsistency between word embedding training and usage space?
 - The objective we optimize during training is not really the one we use
 - Regardless of the different training objective, Word2Vec, GloVe and fastText all optimize the embedding **dot product** during training, but **cosine similarity** is what actually used in applications

Embedding dot product is optimized during training

$$p(w_O|w_I) = \frac{\exp(v'_{w_O}{}^\top v_{w_I})}{\sum_{w=1}^W \exp(v'_w{}^\top v_{w_I})}$$

Word2Vec

$$J = \sum_{i,j=1}^V f(X_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2$$

GloVe

$$s(w, c) = \sum_{g \in \mathcal{G}_w} \mathbf{z}_g^\top \mathbf{v}_c$$

fastText

Motivations

- ❑ What is the consequence of the inconsistency between word embedding training and usage space?
 - ❑ The objective we optimize during training is not really the one we use
 - ❑ E.g. Consider two pairs of words, A: lover-quarrel; B: rock-jazz. Pair B has higher ground truth similarity than pair A in WordSim353 (a benchmark testset)
 - ❑ Word2Vec assigns higher dot product to pair B, but its cosine similarity is still smaller than pair A

	Metrics	A: <i>lover-quarrel</i>		B: <i>rock-jazz</i>
Training	Dot Product	5.284	<	6.287
Usage	Cosine Similarity	0.637	>	0.628

Inconsistency

Motivations

- Apart from the training/usage space inconsistency issue, previous embedding frameworks only leverage **local contexts** to learn word representations
- Local contexts can only partly define word semantics in unsupervised word embedding learning

If I hear someone screwing with my car (ie, setting off the **alarm**) and **taunting** me to come out, you can be very sure that my Colt Delta Elite will also be coming with me. It is not the screwing with the car that would get them **shot**, it is the potential physical **danger**. If they are **taunting** like that, it's very possible that they also intend to **rob** me and or do other physically *harmful* things. Here in Houston last year a woman heard the sound of someone ...


Local contexts of
“harmful”



Introduction

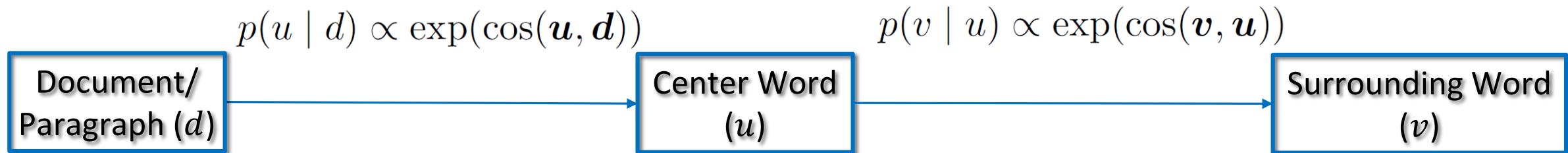
- ❑ In this work, we propose “Spherical Text Embedding”
 - ❑ “**Spherical**”: Embeddings are trained on the unit sphere, where vector norms are ignored and directional similarity is directly optimized
 - ❑ “**Text Embedding**”: Instead of training word embeddings only, we jointly train paragraph (document) embeddings with word embeddings to capture the local and global contexts in text embedding
- ❑ Contributions:
 - ❑ A spherical generative model that jointly exploits word-word (local) and word-paragraph (global) co-occurrence statistics
 - ❑ An efficient optimization algorithm in the spherical space with convergence guarantee
 - ❑ State-of-the-art performances on various text embedding applications

Outline

- ❑ Preliminaries
- ❑ Motivations & Introduction
- ❑ Model: Spherical Text Embedding 
- ❑ Method: Optimization on the Sphere
- ❑ Experiments
- ❑ Conclusions

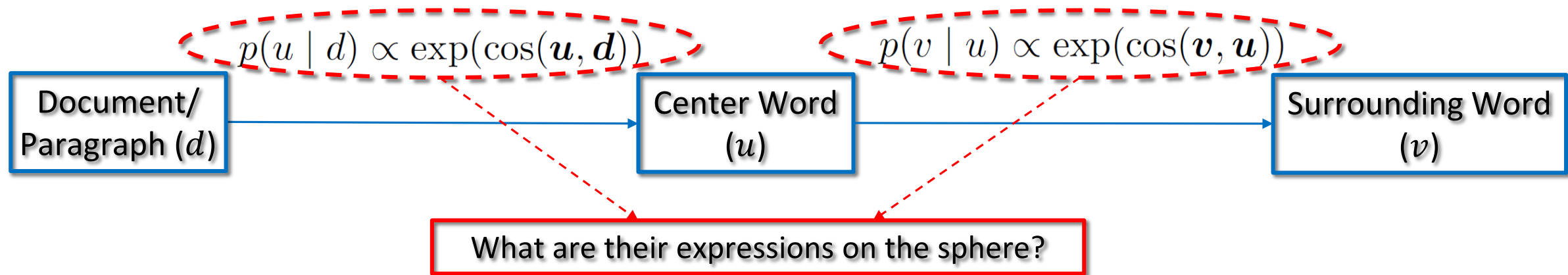
Model: Spherical Text Embedding

- We design a generative model on the sphere that follows how humans write articles:
 - We first have a general idea of the paragraph/document, and then start to write down each word in consistent with not only the paragraph/document, but also the surrounding words
 - Assume a two-step generation process:



Model: Spherical Text Embedding

- How to define the generative model in the spherical space?



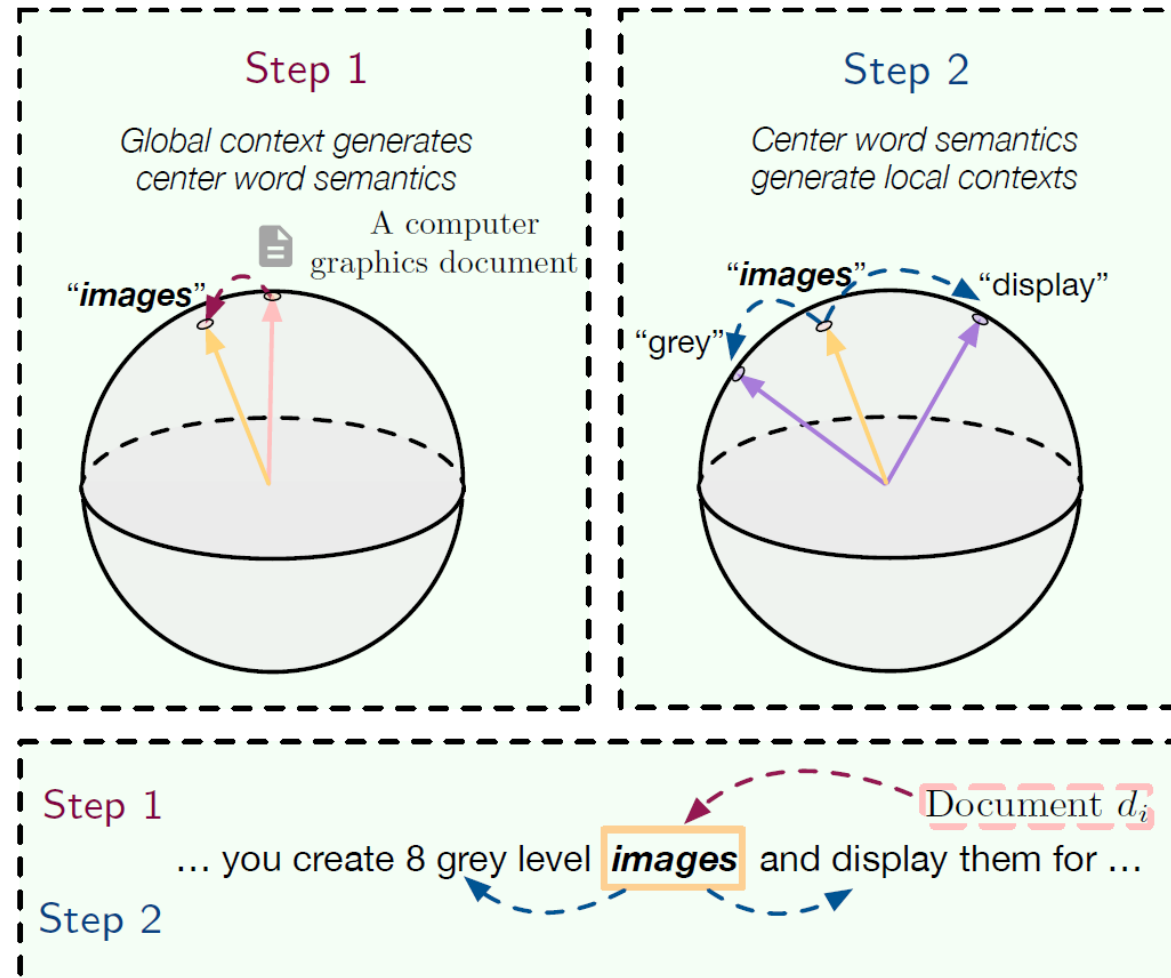
- We prove a theorem connecting the above generative model with a spherical probability distribution:

Theorem 1. When the corpus has infinite vocabulary, *i.e.*, $|V| \rightarrow \infty$, the analytic forms of $p(u | d) \propto \exp(\cos(\mathbf{u}, \mathbf{d}))$ and $p(v | u) \propto \exp(\cos(\mathbf{v}, \mathbf{u}))$ are given by the von Mises-Fisher (vMF) distribution with the prior embedding as the mean direction and constant 1 as the concentration parameter, *i.e.*,

$$\lim_{|V| \rightarrow \infty} p(v | u) = \text{vMF}_p(\mathbf{v}; \mathbf{u}, 1), \quad \lim_{|V| \rightarrow \infty} p(u | d) = \text{vMF}_p(\mathbf{u}; \mathbf{d}, 1).$$

Model: Spherical Text Embedding

□ Understanding the spherical generative model



Model: Spherical Text Embedding

□ Training objective:


- The final generation probability:

$$p(v, u \mid d) = p(v \mid u) \cdot p(u \mid d) = \text{vMF}_p(\mathbf{v}; \mathbf{u}, 1) \cdot \text{vMF}_p(\mathbf{u}; \mathbf{d}, 1)$$

- Maximize the log-probability of a real co-occurred tuple (v, u, d) , while minimize that of a negative sample (v, u', d) , with a max-margin loss:

$$\begin{aligned} \mathcal{L}_{\text{joint}}(\mathbf{u}, \mathbf{v}, \mathbf{d}) &= \max \left(0, m - \underbrace{\log(c_p(1) \exp(\cos(\mathbf{v}, \mathbf{u})) \cdot c_p(1) \exp(\cos(\mathbf{u}, \mathbf{d})))}_{\text{Positive Sample}} \right. \\ &\quad \left. + \underbrace{\log(c_p(1) \exp(\cos(\mathbf{v}, \mathbf{u}')) \cdot c_p(1) \exp(\cos(\mathbf{u}', \mathbf{d})))}_{\text{Negative Sample}} \right) \\ &= \max(0, m - \cos(\mathbf{v}, \mathbf{u}) - \cos(\mathbf{u}, \mathbf{d}) + \cos(\mathbf{v}, \mathbf{u}') + \cos(\mathbf{u}', \mathbf{d})), \end{aligned}$$

Outline

- ❑ Preliminaries
- ❑ Motivations & Introduction
- ❑ Model: Spherical Text Embedding
- ❑ Method: Optimization on the Sphere 
- ❑ Experiments
- ❑ Conclusions

Method: Optimization on the Sphere

- The constrained optimization problem:

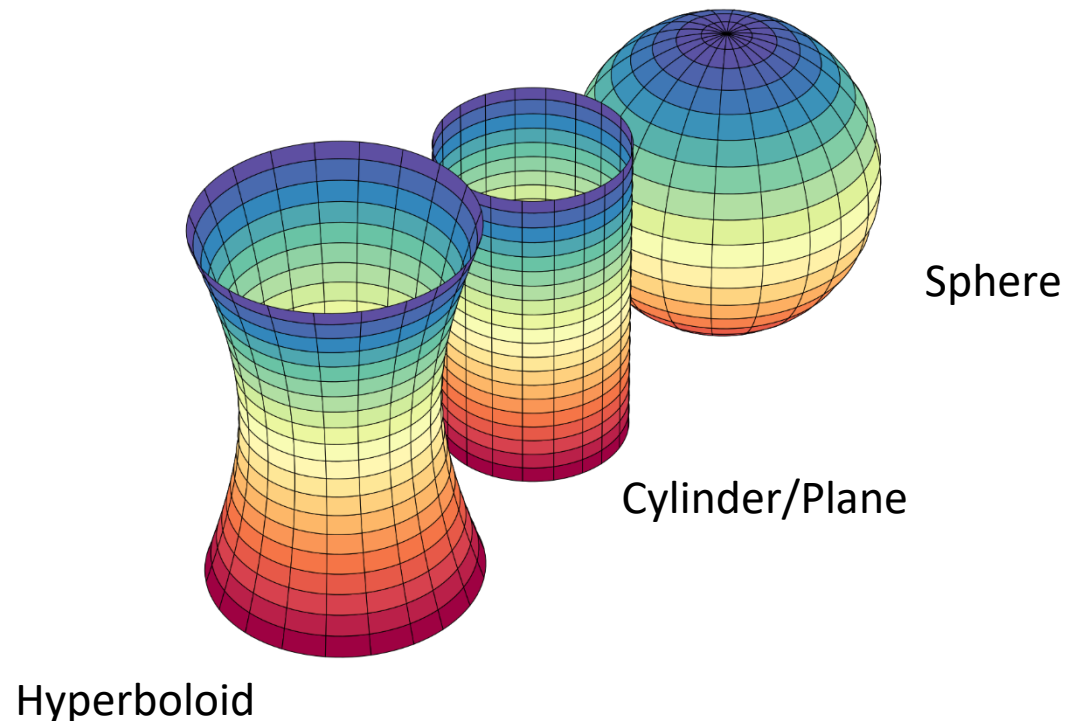
$$\min_{\Theta} \mathcal{L}_{\text{joint}}(\Theta) \quad \text{s.t.} \quad \forall \theta \in \Theta : \|\theta\| = 1 \quad \Theta = \{\mathbf{u}_i\}_{i=1}^{|V|} \cup \{\mathbf{v}_i\}_{i=1}^{|V|} \cup \{\mathbf{d}_i\}_{i=1}^{|\mathcal{D}|}$$

- Challenge: Parameters must be always updated on the sphere, but Euclidean optimization methods (e.g. SGD) are not constrained on a curvature space
- Need to consider the nature of the spherical space

Method: Optimization on the Sphere

- Riemannian manifold:

- The sphere is a Riemannian manifold with constant positive curvature



Method: Optimization on the Sphere

□ Riemannian optimization with Riemannian SGD:

□ Riemannian gradient:

$$\text{grad } f(\mathbf{x}) := (I - \mathbf{x}\mathbf{x}^\top) \nabla f(\mathbf{x})$$

□ Exponential mapping (maps from the tangent plane to the sphere):

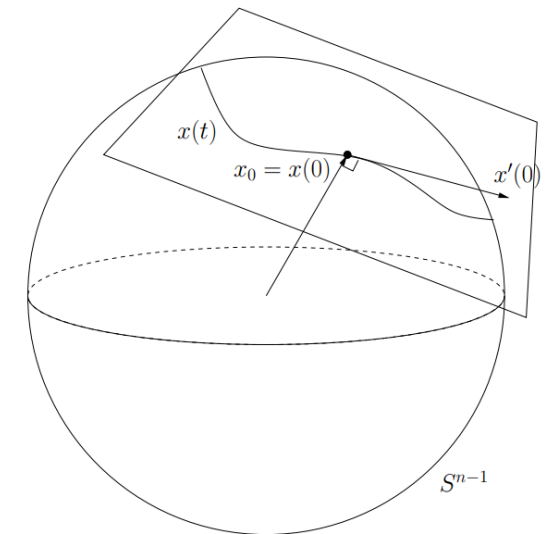
$$\exp_{\mathbf{x}}(\mathbf{z}) := \begin{cases} \cos(\|\mathbf{z}\|)\mathbf{x} + \sin(\|\mathbf{z}\|)\frac{\mathbf{z}}{\|\mathbf{z}\|}, & \mathbf{z} \in T_{\mathbf{x}}\mathbb{S}^{p-1} \setminus \{\mathbf{0}\}, \\ \mathbf{x}, & \mathbf{z} = \mathbf{0}. \end{cases}$$

□ Riemannian SGD:

$$\mathbf{x}_{t+1} = \exp_{\mathbf{x}_t}(-\eta_t \text{grad } f(\mathbf{x}_t))$$

□ Retraction (first-order approximation of the exponential mapping):

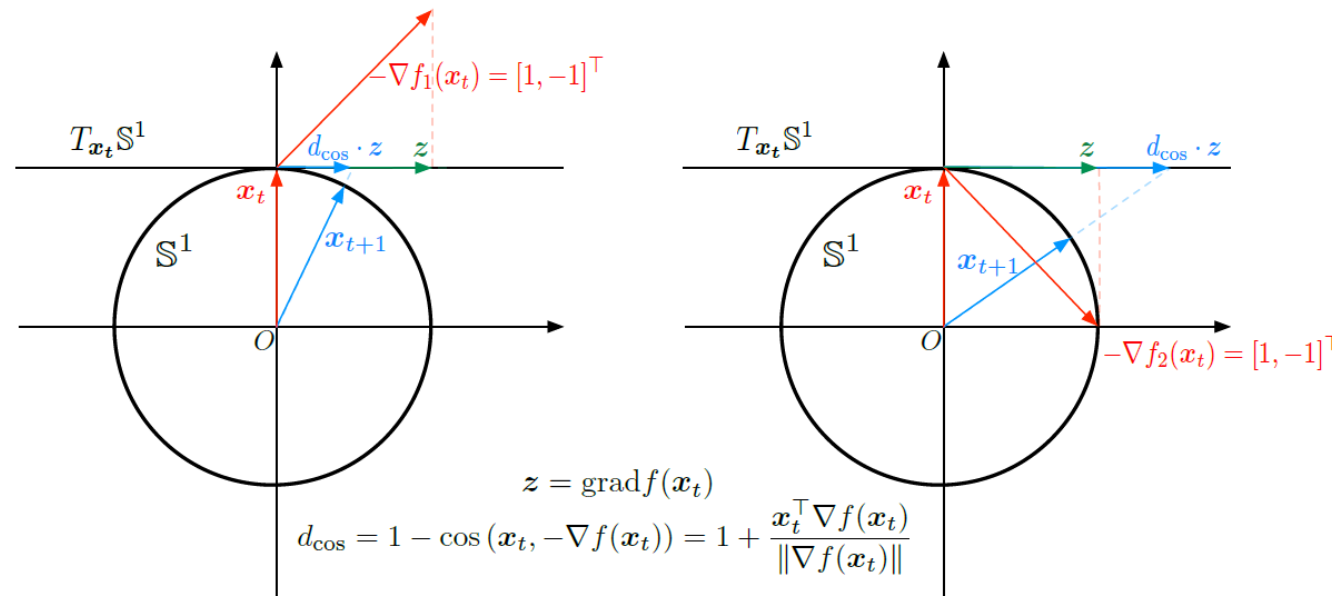
$$R_{\mathbf{x}}(\mathbf{z}) := \frac{\mathbf{x} + \mathbf{z}}{\|\mathbf{x} + \mathbf{z}\|}$$



Method: Optimization on the Sphere

□ Training details:

- Incorporate angular distances into Riemannian optimization



- Multiply the Euclidean gradient with its angular distance from the current point

$$x_{t+1} = R_{x_t} \left(-\eta_t \left(1 + \frac{x_t^\top \nabla f(x_t)}{\|\nabla f(x_t)\|} \right) (I - x_t x_t^\top) \nabla f(x_t) \right).$$

Method: Optimization on the Sphere


- Convergence guarantee of the optimization procedure:

Theorem 2. When the update rule given by Equation (7) is applied to $\mathcal{L}(x)$, and the learning rate satisfies the usual condition in stochastic approximation, *i.e.*, $\sum_t \eta_t^2 < \infty$ and $\sum_t \eta_t = \infty$, x converges almost surely to a critical point x^* and $\text{grad } \mathcal{L}(x)$ converges almost surely to 0, *i.e.*,

$$\Pr \left(\lim_{t \rightarrow \infty} \mathcal{L}(x_t) = \mathcal{L}(x^*) \right) = 1, \quad \Pr \left(\lim_{t \rightarrow \infty} \text{grad } \mathcal{L}(x_t) = 0 \right) = 1.$$



Outline

- ❑ Preliminaries
- ❑ Motivations & Introduction
- ❑ Model: Spherical Text Embedding
- ❑ Method: Optimization on the Sphere
- ❑ Experiments 
- ❑ Conclusions



Experiments

- ❑ Word similarity:
 - ❑ Train 100-d word embedding on the latest Wikipedia dump (~13G)
 - ❑ Compute embedding **cosine similarity** between word pairs to obtain a ranking of similarity
 - ❑ Benchmark datasets contain human rated similarity scores
 - ❑ The more similar the two rankings are, the better embedding reflects human thoughts
 - ❑ **Spearman's rank correlation** is used to measure the ranking similarity

Experiments

❑ Word similarity baselines:

❑ Euclidian Space:

- ❑ Word2Vec: Distributed representations of words and phrases and their compositionality. In NIPS, 2013
- ❑ GloVe: Glove: Global vectors for word representation. In EMNLP, 2014
- ❑ fastText: Enriching word vectors with subword information. TACL, 2017
- ❑ BERT: Pre-training of deep bidirectional transformers for language understanding. In NAACL, 2019

❑ Poincare Space:

- ❑ Poincaré glove: Hyperbolic word embeddings. In ICLR, 2019

❑ Spherical Space:

- ❑ JoSE (our full model)

Experiments

□ Word similarity results:

Table 1: Spearman rank correlation on word similarity evaluation.

Embedding Space	Model	WordSim353	MEN	SimLex999
Euclidean	Word2Vec	0.711	0.726	0.311
	GloVe	0.598	0.690	0.321
	fastText	0.697	0.722	0.303
	BERT	0.477	0.594	0.287
Poincaré	Poincaré GloVe	0.623	0.652	0.321
Spherical	JoSE	0.739	0.748	0.339

□ Why does BERT fall behind on this task?

- BERT learns contextualized representations, but word similarity is conducted in a context-free manner
- BERT is optimized on specific downstream tasks like predicting masked words and sentence relationships, which have no direct relation to word similarity



Experiments

❑ Document Clustering:

- ❑ Train document embedding on 20News dataset (20 classes)
- ❑ Perform K-means and Spherical K-means (SK-means)
- ❑ Metrics: Mutual Information (MI), Normalized Mutual Information (NMI), Adjusted Rand Index (ARI), and Purity
- ❑ Run clustering 10 times and report the above metrics with mean and standard deviation



Experiments

❑ Document clustering baselines:

❑ Euclidian Space:

- ❑ Avg. Word2Vec: Use the averaged word embedding of Word2Vec as document embedding
- ❑ SIF: Simple but tough-to-beat baseline for sentence embeddings. In ICLR, 2017
- ❑ BERT: Pre-training of deep bidirectional transformers for language understanding. In NAACL, 2019
- ❑ Doc2Vec: Distributed representations of sentences and documents. In ICML, 2014

❑ Spherical Space:

- ❑ JoSE (our full model)

Experiments

Document clustering results:

Table 2: Document clustering evaluation on the 20 Newsgroup dataset.

Embedding	Clus. Alg.	MI	NMI	ARI	Purity
Avg. W2V	K-Means	1.299 ± 0.031	0.445 ± 0.009	0.247 ± 0.008	0.408 ± 0.014
	SK-Means	1.328 ± 0.024	0.453 ± 0.009	0.250 ± 0.008	0.419 ± 0.012
SIF	K-Means	0.893 ± 0.028	0.308 ± 0.009	0.137 ± 0.006	0.285 ± 0.011
	SK-Means	0.958 ± 0.012	0.322 ± 0.004	0.164 ± 0.004	0.331 ± 0.005
BERT	K-Means	0.719 ± 0.013	0.248 ± 0.004	0.100 ± 0.003	0.233 ± 0.005
	SK-Means	0.854 ± 0.022	0.289 ± 0.008	0.127 ± 0.003	0.281 ± 0.010
Doc2Vec	K-Means	1.856 ± 0.020	0.626 ± 0.006	0.469 ± 0.015	0.640 ± 0.016
	SK-Means	1.876 ± 0.020	0.630 ± 0.007	0.494 ± 0.012	0.648 ± 0.017
JoSE	K-Means	1.975 ± 0.026	0.663 ± 0.008	0.556 ± 0.018	0.711 ± 0.020
	SK-Means	1.982 ± 0.034	0.664 ± 0.010	0.568 ± 0.020	0.721 ± 0.029

- Embedding quality is generally more important than clustering algorithms:
 - Using spherical K-Means only gives marginal performance boost over K-Means
 - JoSE embedding remains optimal regardless of clustering algorithms

Experiments

❑ Document Classification:

- ❑ Train document embedding on 20News (20 classes) and Movie review (2 classes) dataset
- ❑ Perform k-NN classification ($k=3$)
- ❑ Metrics: Macro-F1 & Micro-F1
- ❑ Baselines: Same with document clustering

Experiments

□ Document classification results:

Table 3: Document classification evaluation using k -NN ($k = 3$).

Embedding	20 Newsgroup		Movie Review	
	Macro-F1	Micro-F1	Macro-F1	Micro-F1
Avg. W2V	0.630	0.631	0.712	0.713
SIF	0.552	0.549	0.650	0.656
BERT	0.380	0.371	0.664	0.665
Doc2Vec	0.648	0.645	0.674	0.678
JoSE	0.703	0.707	0.764	0.765

- We observe similar comparison results with k from $[1, 10]$
- k -NN is non-parametric and directly reflect how well the topology of the embedding space captures document-level semantics

Experiments

❑ Training efficiency:

Table 4: Training time (per iteration) on the latest Wikipedia dump.


Word2Vec	GloVe	fastText	BERT	Poincaré GloVe	JoSE
0.81 hrs	0.85 hrs	2.11 hrs	> 5 days	1.25 hrs	0.73 hrs

❑ Why is JoSE training efficient?

- ❑ Other models' objectives contain many non-linear operations (Word2Vec & fastText's objectives involve exponential functions; GloVe's objective involves logarithm functions), while JoSE only has linear terms in the objective



Outline

- ❑ Preliminaries
- ❑ Motivations & Introduction
- ❑ Model: Spherical Text Embedding
- ❑ Method: Optimization on the Sphere
- ❑ Experiments
- ❑ Conclusions 

Conclusions

- ❑ In this work, we introduce a spherical text embedding framework
 - ❑ Address the discrepancy between the training procedure and the practical usage of text embedding
 - ❑ Introduce a spherical generative model to jointly learn word and paragraph embedding
 - ❑ Develop an efficient Riemannian optimization method to train text embedding on the unit hypersphere
 - ❑ Achieves state-of-the-art performances on common text embedding applications
- ❑ Future work:
 - ❑ Exploit spherical embedding space for other tasks like lexical entailment
 - ❑ Incorporate other signals such as subword information into spherical text embedding
 - ❑ Benefit other supervised tasks: Word embedding is commonly used as the first layer in DNN
 - ❑ Add norm constraints to word embedding layer
 - ❑ Apply Riemannian optimization when fine-tuning the word embedding layer



Thanks!