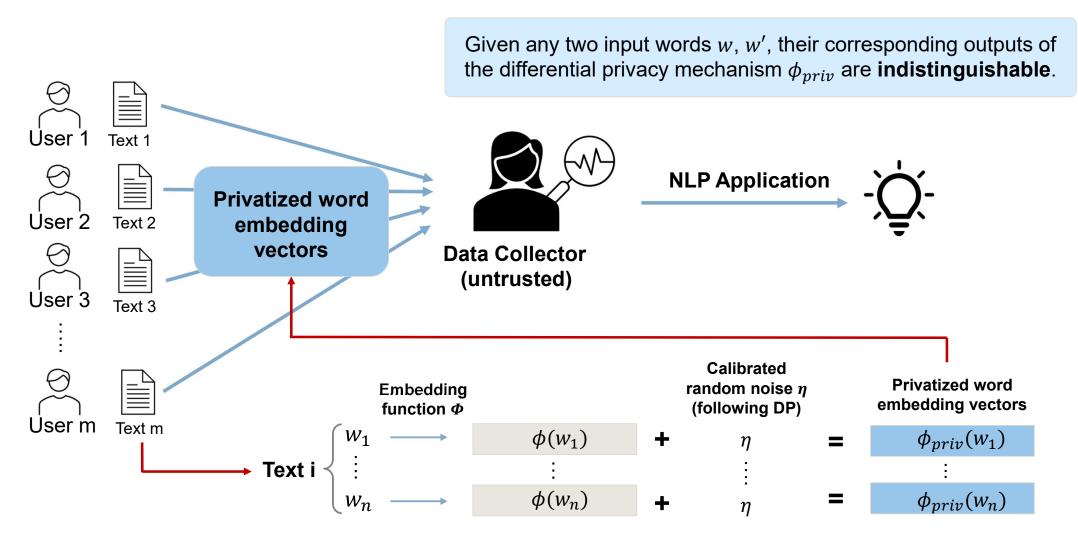




- 1. Motivation
- 2. Research Questions
- 3. Methodology
- 4. Initial Findings
- 5. Next Steps
- 6. Timeline

Motivation – Differential Privacy for Word Embedding Vectors (1)





Motivation - Differential Privacy for Word Embedding Vectors (2)



 $\Phi(w)$

(Embedding Function)

+

η

(Calibrated Random Noise)

 $\Phi_{priv}(w)$

(Randomized Mechanism)

Example: Laplace Mechanism

An embedding function ${m \Phi}({m w})$ with sensitivity ${m \Delta}_{m \phi}$

 $\eta \sim Lap$

 Φ_{priv} is ϵ -differentially private

Sensitivity describes the maximum possible change in a function's output when the input is changed.

Variance depends on sensitivity Δ_{ϕ} and privacy budget ε

The probability of a specific output on any two inputs w, w' can differ by at most a **multiplicative factor** dependent on ε .

Motivation - Differential Privacy for Word Embedding Vectors (3)



 $\Phi(w)$

(Embedding Function)



η

(Calibrated Random Noise)



 $\Phi_{priv}(w)$

(Randomized Mechanism)



Large sensitivity Δ_{ϕ} leads to large noise η



Bound sensitivity



How can sensitivity be bounded?

How does the bounding affect the Privacy-Utility Trade-off?



 $\Phi_{priv}(w)$ is not a "real" word



Map $\Phi_{priv}(w)$ to similar embedding vector associated with a "real" word



- What can the mapping look like?
 - How does this mapping affect the Privacy-Utility Trade-off?



- 1. Motivation
- 2. Research Questions
- 3. Methodology
- 4. Initial Findings
- 5. Next Steps
- 6. Timeline

Research Questions



RQ1	What approaches are there to privatize word embeddings by perturbing word vector representations?
RQ2	How can we make these privatized word embeddings more effective?
RQ3	What is the effect of different approaches to estimating sensitivity on privacy and utility for downstream NLP tasks?
RQ4	What are the implications on privacy and utility resulting from mapping noisy word embeddings to similar embedding vectors which are associated with real words?



- 1. Motivation
- 2. Research Questions
- 3. Methodology
- 4. Initial Findings
- 5. Next Steps
- 6. Timeline

Methodology – Literature Review



Domain:

- Text data
- Image data
- Location data

Differential Privacy Method:

- Embedding perturbation
- Gradient perturbation



Goals

RQ1

Approaches to privatize word embeddings by **perturbing** word vector representations

RQ2

Ideas to increase the utility
of the privatized word
embeddings
(focus on sensitivity- and
mapping-related ideas)

Methodology - Experiments



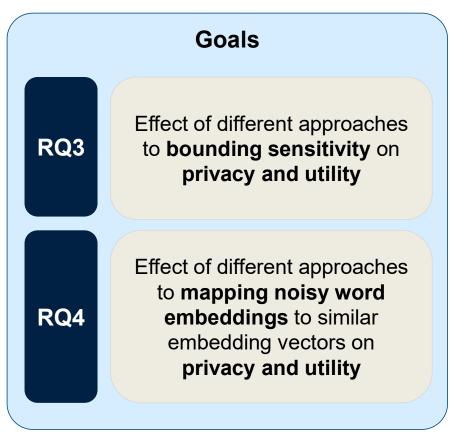
Empirically compare privatized embedding vectors

Utility Evaluation

- Sentiment Analysis
- Topic Classification

Privacy Evaluation

- Identification of authors' gender and age
- Identification of named entities





- 1. Motivation
- 2. Research Questions
- 3. Methodology
- 4. Initial Findings
- 5. Next Steps
- 6. Timeline

Initial Findings – Bounding Sensitivity



Concatenate Φ with a function f_{bound}

 $\Phi(w)$ has unbounded sensitivity



$$\widetilde{\Phi}(w) = f_{bound} \circ \Phi(w) \text{ s.th. } \widetilde{\Phi}(w) \in [a, b]^d$$

s.th.
$$\widetilde{\Phi}(w) \in [a,b]^d$$

 $\widetilde{\Phi}(w)$ has bounded sensitivity

Bounding so	ensitivity using $f_{\it bound}$	Sensitivity	
Bound the	By observed range	$\Delta_{\widetilde{\phi}} \leq d*2*v_{max},$ where v_{max} is the maximum observed value	
embedding space	Normalize to $[0,1]^d$	$\Delta_{\widetilde{\phi}} \leq d$	
	Normalize to unit length	$\Delta_{\widetilde{\Phi}} \leq 2$	
Dimensionality Reduction	Johnson-Lindenstrauss Lemma	$\Delta_{\widetilde{\Phi}} \leq d'*(b-a),$ where $d' < d$	

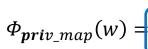
$$f_{bound}(y) = \frac{y - \min(y)}{\max(y) - \min(y)}$$

Initial Findings – Mapping to "real" words



 $\Phi_{priv}(w)$ does not correspond to a word embedding vector





 $\Phi_{priv_map}(w) = f_{map} \circ \Phi_{priv}(w), \text{ s.th } , \Phi_{priv_map}(w) \in Im(\Phi)$

Concatenate Φ_{priv} with a function f_{map}

 $\Phi_{priv_map}(w)$ corresponds to a "real" word's embedding vector

Mapping	Resulting embedding vector
Map to nearest neighbor embedding	Similar embedding vector associated with "real" word
Randomly output first or second nearest neighbor	Similar embedding vector associated with "real" word
No mapping	Noisy embedding vector



- 1. Motivation
- 2. Research Questions
- 3. Methodology
- 4. Initial Findings
- 5. Next Steps
- 6. Timeline

Next Steps - Experiments

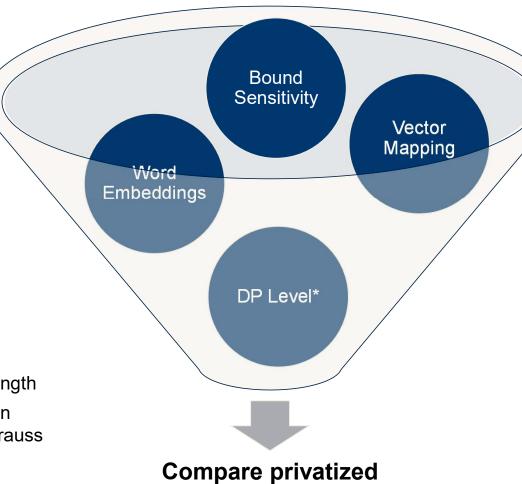


Word Embeddings:

- GloVe embeddings
- FastText embeddings

Bound Sensitivity:

- observed range
- normalize to $[0,1]^d$
- normalize to unit length
- dimension reduction (Johnson-Lindenstrauss Lemma)



embedding vectors

Vector Mapping:

- map to nearest neighbor embedding
- randomly output first or second nearest neighbor
- no vector mapping

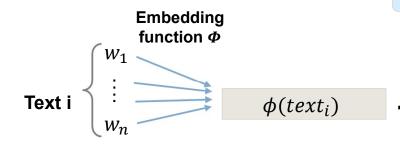
DP Level:

- Word level privacy
- Document/User level privacy

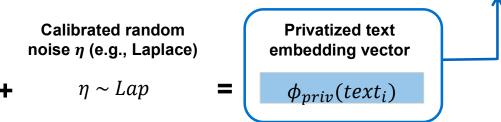
*Differential Privacy Levels – Document vs. Word Level Privacy



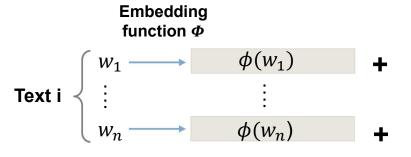
Document level privacy



Given any two input texts t, t', their corresponding outputs of the differential privacy mechanism ϕ_{priv} are ε -indistinguishable.



Word level privacy



Calibrated random noise η (e.g., Laplace)

$$\eta \sim Lap$$

$$\vdots$$

$$\eta \sim Lap$$

Privatized word embedding vectors

$$\phi_{priv}(w_1)$$
 \vdots
 $\phi_{priv}(w_n)$

Document level privacy

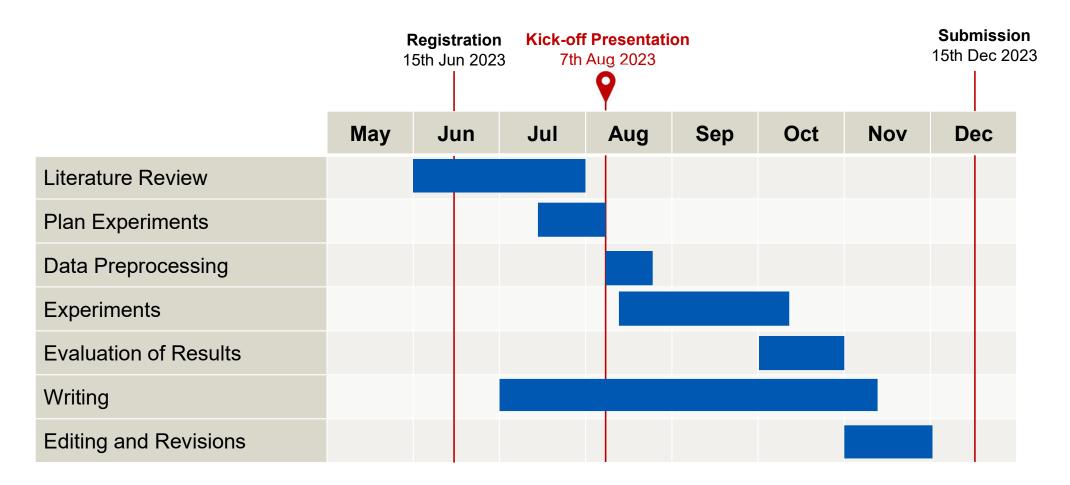
Given any two input texts t, t', their corresponding outputs of the differential privacy mechanism ϕ_{priv} are $\varepsilon * n$ -indistinguishable.



- 1. Motivation
- 2. Research Questions
- 3. Methodology
- 4. Initial Findings
- 5. Next Steps
- 6. Timeline

Timeline







Motivation - Differential Privacy for Word Embedding Vectors



$$\Phi(w)$$

+

η

 $\Phi_{priv}(w)$

(Embedding Function)

(Calibrated Random Noise)

(Randomized Mechanism)

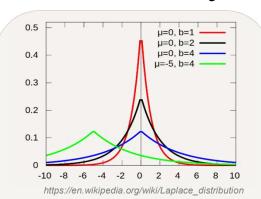
Example: Laplace Mechanism

An embedding function $\Phi(w)$ with sensitivity Δ_{ϕ} :

$$\Delta_{\phi} = \max_{x, x' \in X} \lVert \Phi(w) - \Phi(w') \rVert_{1}$$

Sensitivity describes the maximum possible change in a function's output, when the input is changed.

$$\eta \sim Lap(\mu, b),$$
 where $\mu = 0$, $\mathbf{b} = \frac{\Delta_{\phi}}{\varepsilon}$



Variance depends on sensitivity Δ_{ϕ} and privacy budget ε

 $Φ_{priv}$ is ε-differentially private if for all pairs of inputs w, w' and all possible outputs z:

$$\frac{P[\Phi_{priv}(w) = z]}{P[\Phi_{priv}(w') = z]} \le \exp(\varepsilon)$$

The probability of a specific output on any two inputs w, w' can differ by at most a **multiplicative factor of** $exp(\varepsilon)$.

Approaches for Bounding Sensitivity



© sebis

 $\Phi(w)$ has unbounded sensitivity



$$\widetilde{\Phi}(w) = f_{bound} \circ \Phi(w)$$
, s.th. $\widetilde{\Phi}(w) \in [a, b]^d$

Approaches to bounding sensitivity		Function f_{bound}	Sensitivity $arDelta_{\widetilde{\phi}}$
Bound the	By observed range	Define v_{max} as the maximum observed value. $f_{bound}(y)_i = \begin{cases} y_i, & y_i \in [-v_{max}, v_{max}] \\ -v_{max}, & y_i < -v_{max} \\ v_{max}, & y_i > v_{max} \end{cases}$	$\Delta_{\widetilde{\Phi}} \le d * 2 * v_{max}$
embedding space	Normalize to $[0,1]^d$	$f_{bound}(y) = \frac{y - \min(y)}{\max(y) - \min(y)}$	$\Delta_{\widetilde{\Phi}} \le d$
	Normalize to unit length	$f_{bound}(y) = \frac{y}{\ y\ }$	$\Delta_{\widetilde{\Phi}} \leq 2$
Dimensionality Reduction	Johnson- Lindenstrauss Lemma	$f_{bound}(y) = A * x$, where $A \in \mathbb{R}^{k \times m}$ with a_{ij} drawn from $\left\{-\frac{1}{\sqrt{k}}, +\frac{1}{\sqrt{k}}\right\}$ iid. and $k = \left[\frac{8 \ln(n)}{\varepsilon^2 - 2\frac{\varepsilon^3}{3}}\right]$	$\Delta_{\widetilde{\Phi}} \leq d'*(b-a),$ where $d' < d$

Approaches for Mapping to "real" words



 $\Phi_{priv}(w)$ does not correspond to a word embedding vector

 $\Phi_{\textit{priv}_\textit{map}}(w) = f_{\textit{map}} \circ \Phi_{\textit{priv}}(w), \, \text{s.th.}, \\ \Phi_{\textit{priv}_\textit{map}}(w) \in Im(\Phi)$

Mapping	Function f_{map}	Embedding Vector
No mapping	$f_{map}(y) = y$	Noisy embedding vector
Map to nearest neighbor embedding	$f_{map}(y) = \underset{w \in \mathcal{W}\setminus\{x\}}{\operatorname{argmin}} \ y - \Phi(w)\ $	Similar embedding vector associated with "real" word
Randomly output first or second nearest neighbor	$\begin{split} f_{map}(y) &= \begin{cases} w_1, & \textit{with prob.p} \\ w_2, & \textit{with prob.} \ 1-p \end{cases} \\ \text{where } w_1 &= \underset{w \in \mathcal{W} \setminus \{x\}}{\operatorname{argmin}} \parallel y - \Phi(w) \parallel \\ w_2 &= \underset{w \in \mathcal{W} \setminus \{x, w_1\}}{\operatorname{argmin}} \parallel y - \Phi(w) \parallel \end{split}$	Similar embedding vector associated with "real" word

Approaches for Bounding Sensitivity

 $\Phi(w)$ has unbounded sensitivity



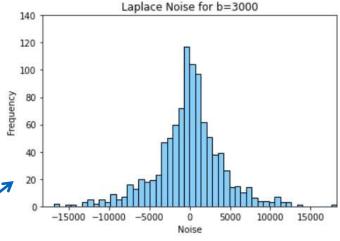
$$\widetilde{\Phi}(w) = f_{bound} \circ \Phi(w)$$
, s.th. $\widetilde{\Phi}(w) \in [a, b]^d$

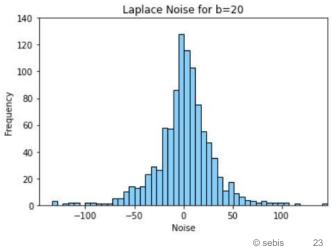
Privacy budget: $\varepsilon = 0.1$

Embedding function: $\Phi(w)$: $\mathcal{W} \to \mathbb{R}^{300}$

Noise: $\eta \sim Lap(\mu, b)$, where $\mu = 0$, $b = \frac{\Delta_{\phi}}{c}$

Approaches to bounding sensitivity	Sensitivity $arDelta_{\widetilde{\phi}}$	Variance of Laplace noise
Normalize to $[0,1]^d$	$\Delta_{\widetilde{\phi}} \leq d$	$rac{arDelta_{\phi}}{arepsilon}=3000$
Normalize to unit length	$\Delta_{\widetilde{\Phi}} \leq 2$	$rac{arDeta_{m{\phi}}}{arepsilon}=20$





Johnson-Lindenstrauss Lemma



Goal: Bound the worst-case distortion of distances during dimensionality reduction

By the Johnson-Lindenstrauss Lemma:

- $P = \{x_1, x_2, ..., x_n\}$ a set of n points in \mathbb{R}^m
- $\xi \in (0,1)$
- $k = \left[\frac{8 \ln(n)}{\xi^2 2\frac{\xi^3}{3}} \right]$

There exists a randomly generated mapping $f: \mathbb{R}^m \to \mathbb{R}^k$ such that with probability at least $1 - \frac{1}{n^2}$:

$$(1-\xi)\|x_i-x_j\|_2 \le \|f(x_i)-f(x_j)\|_2 \le (1+\xi)\|x_i-x_j\|_2$$

for all i, j = 1, 2, ..., n.

Possible choice for mapping f:

 $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$, where A is a k x m matrix with a_{ij} iid. from $\left\{-\frac{1}{\sqrt{k}}, +\frac{1}{\sqrt{k}}\right\}$

Unbounded Sensitivity



$$\Phi(w)$$



η

 $\Phi_{priv}(w)$

(Randomized Mechanism)

(Embedding Function)

(Calibrated Random Noise)

Sensitivity: $\Delta_{\phi} = \max_{w.w' \in \mathcal{W}} \lVert \Phi(w) - \Phi(w') \rVert$

Problem: Unbounded Sensitivity

$$\Phi(w) \in (-\infty, +\infty)^d \Rightarrow \Delta_{\phi} = \infty$$



Approach: Bound Sensitivity

$$\widetilde{\Phi}(w) = f_{bound} \circ \Phi(w), \text{ s.th. } \widetilde{\Phi}(w) \in [a,b]^d \quad \Rightarrow \quad \Delta_{\widetilde{\Phi}} = d * (b-a) < \infty$$

$\Phi_{priv}(w)$ is not a "real" word



$$\Phi(w)$$

(Embedding Function)

+

η

(Calibrated Random Noise)

 $\Phi_{priv}(w)$

(Randomized Mechanism)

Problem: $\Phi_{priv}(w)$ is not a "real" word

 $\Phi_{priv_map}(w) \notin Im(\Phi)$



Approach: Map $\Phi_{priv}(w)$ to similar embedding vectors associated with real words

 $\Phi_{priv_map}(w) = f_{map} \circ \Phi_{priv}(w), \text{ s.th.}, \Phi_{priv_map}(w) \in Im(\Phi)$