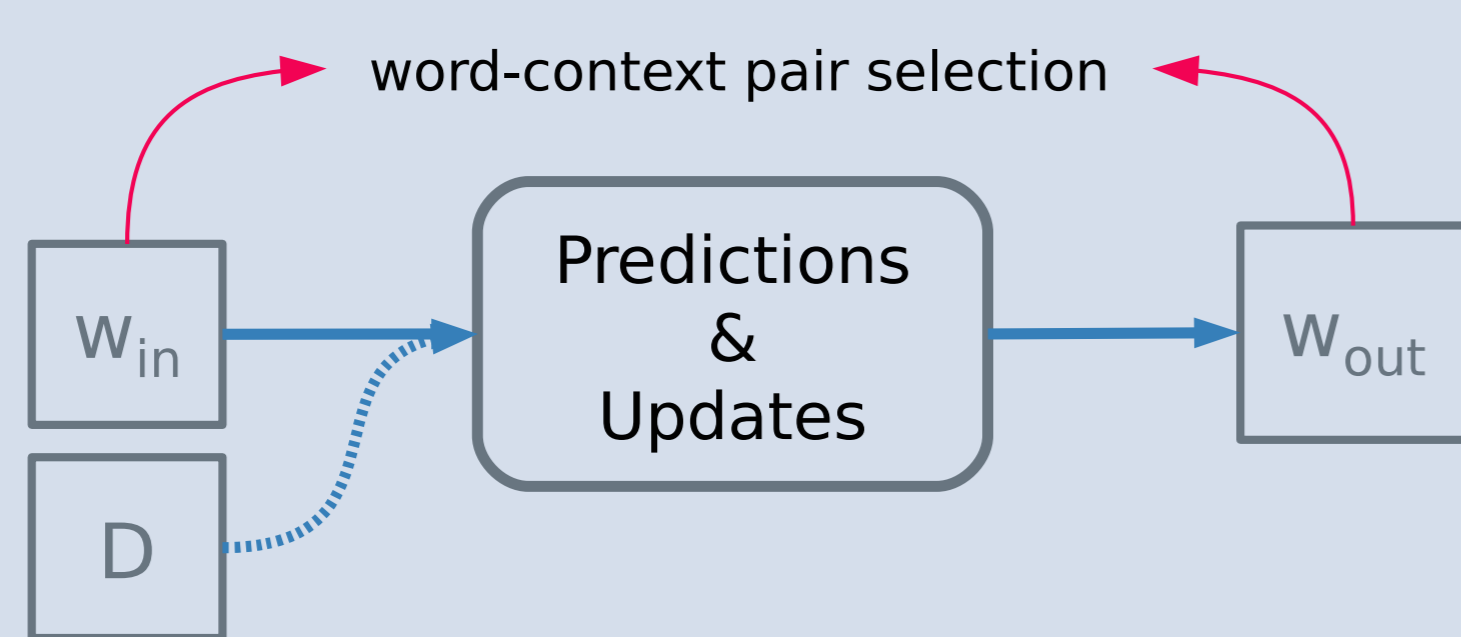


Why Document Embeddings?

- ▶ Embedding models are used widely for learning word representations from vast amounts of **unlabeled** text.
- ▶ Represent meaning of longer pieces of text → embedding composition.
- ▶ Averaging / Syntax-aided composition:
 - ▷ Can work for phrases or short sentences.
 - ▷ Severe loss of semantic information as sequence length increases.
- ▶ CNNs / Recurrent NNs / Hierarchical NNs:
 - ▷ State-of-the-art in many **supervised** tasks.
 - ▷ Computationally **demanding** (need GPUs).
- ▶ Middle ground → Paragraph Vector (Le and Mikolov, 2014).

Paragraph Vector



- ▶ Extension of Skip-gram (Mikolov et al., 2013).
- ▶ Word and document embeddings are learned **jointly** w/o supervision.
- ▶ Words are paired with their **window-based** contexts.
- ▶ Document embeddings are also used to predict each word they contain.
- ▶ **Issue:** Are selected word-context pairs representative of content?

PV
The Ministry of industry is expected to give Elf **Aquitaine**, the oil group controlled by the French state, total clearance to take over the UK petrol-station.

Window-based context

- ▶ Disregards word importance.
- ▶ Implicitly forces doc. embeddings towards frequent words.

This work: Context Sampling Framework

- ▶ We introduce arbitrary contexts via **Context Sampling**.
- ▶ Different sampling **policies** will result in different embedding spaces.

DE_{idf}
The Ministry of industry is expected to give Elf **Aquitaine**, the oil group controlled by the French state, total clearance to take over the UK petrol-station.

IDF Sampling

- ▶ Context words sampled from document-wide *tf.idf* distribution.
- ▶ Doc. embeddings are positioned closer to **content-heavy** words.

DE_{nn}
The Ministry of industry is expected to give Elf **Aquitaine**, the oil group controlled by the French state, total clearance to take over the UK petrol-station.

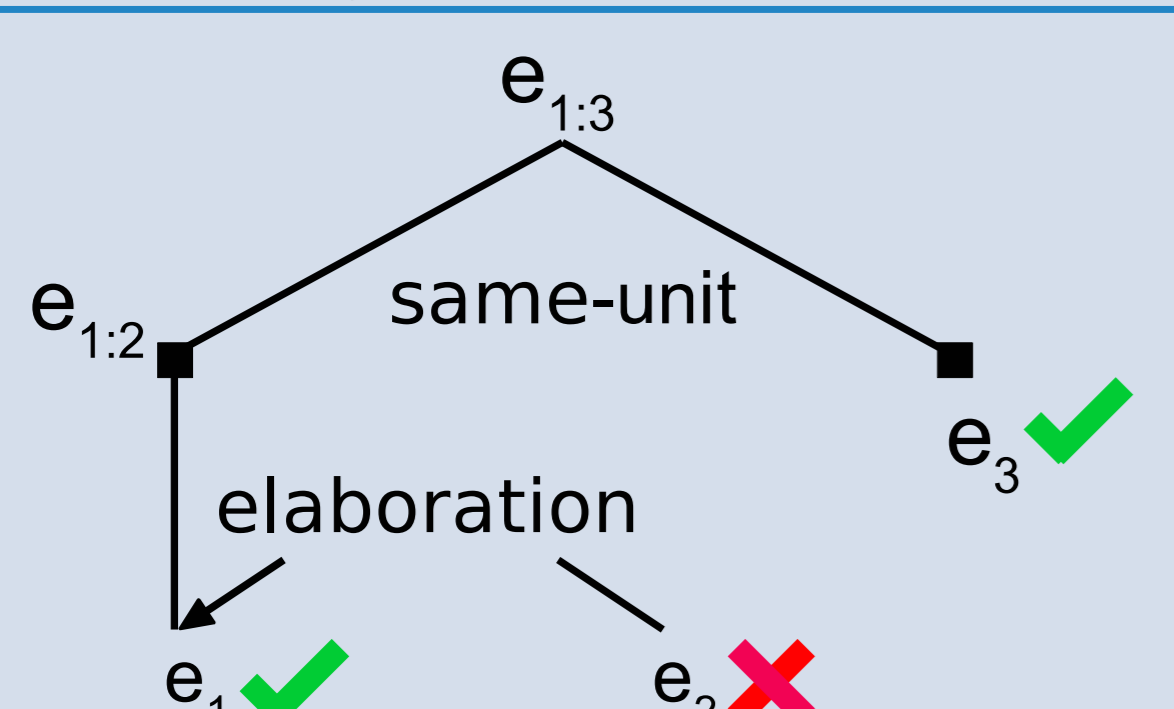
Neighborhood Sampling

- ▶ Incorporates **clustering hypothesis** to context selection.
- ▶ Context words sampled from fixed-size neighborhood of similar documents.
- ▶ Words that do not appear in current document may be used as well.

DE_{disc}
[The Ministry of industry is expected to give Elf **Aquitaine**], e₁ [the oil group controlled by the French state], e₂ [total clearance to take over the UK petrol-station]. e₃

Discourse-based Sampling

- ▶ Attempts to **inject discourse-level** linguistic information.
- ▶ Not all parts of a document are equally important.
- ▶ We **parse** documents using an RST-style discourse parser (Feng and Hirst, 2012).
- ▶ Potentially insignificant **elementary discourse unit** (EDU) types are filtered-out before context selection.



Evaluation

- ▶ Compare the quality of document embeddings learned using PV's window-based contexts and our context sampling policies (DE).
- ▶ We chose 2 **document-centric** tasks:
 - ▷ Ad-Hoc Search
 - ▷ Document Classification

Results: Ad-hoc Search

- ▶ Used 2 established TREC collection for Information Retrieval.
- ▶ Learned **query** and **document** embeddings (no supervision).
- ▶ Documents ranked by their cosine distance to a query in **embedding space**.

Methods	ROBUST			AP88-89		
	MAP	%-change vs. PV	%-change vs. DE _{idf}	MAP	%-change vs. PV	%-change vs. DE _{idf}
PV	0.1179	—	—	0.0938	—	—
DE _{idf}	0.1328	12.6*	—	0.1154	23.0*	—
DE _{q.nn}	0.1693	43.6*	27.5*	0.1442	53.7*	24.9*
DE _{f.nn}	0.1823	54.6*	37.3*	0.1631	73.8*	41.3*

* indicates significant improvement based on a two-tailed t-test with $p < 0.01$.

- ▶ Can provide complementary signal to term-based IR methods.

Results: Document Classification

- ▶ Sentiment Analysis (IMDB) & Topic Classification (RCV1).
- ▶ Learned document embeddings (no supervision).
- ▶ Trained logit classifier using document embeddings as **features**.

Classification Performance

Methods	Accuracy (%)	
	IMDB	RCV1
N-gram	86.52	85.12
RNN-LM	86.61	85.08
PV	88.93	86.95
DE _{idf}	89.29 ^a	87.71 ^a
DE _{nn}	89.34^a	87.75 ^a
DE _{disc.pv}	88.37	87.05 ^a
DE _{disc.idf}	88.82	87.97 ^{abc}
DE _{disc.nn}	88.87	88.01^{abc}

Markers *a*, *b* and *c* denote significant improvements over PV, DE_{idf} and DE_{nn} resp. (one-tail t-test with $p < 0.01$).

Qualitative Evaluation

	PV	DE _{idf}	DE _{disc.idf}
health	<i>care</i>	<i>medical</i>	<i>medical</i>
	<i>medical</i>	<i>physician</i>	<i>hospital</i>
	<i>education</i>	<i>hospital</i>	<i>nhs</i>
politics	<i>benefits</i>	<i>therapy</i>	<i>nursing</i>
	<i>political</i>	<i>political</i>	<i>political</i>
	<i>politicians</i>	<i>party</i>	<i>party</i>
weather	<i>candidature</i>	<i>polls</i>	<i>election</i>
	<i>dirty</i>	<i>election</i>	<i>leader</i>
	<i>cooler</i>	<i>temperatures</i>	<i>temperatures</i>
	<i>warm</i>	<i>meteorologist</i>	<i>dry</i>
	<i>dry</i>	<i>warm</i>	<i>meteorologist</i>
	<i>warmer</i>	<i>rain</i>	<i>precipitation</i>

Based on cosine similarity

- ▶ Ranked words against common RCV1 topics.
- ▶ PV produces embeddings that reflect **co-occurrence** patterns.
- ▶ Document-wide context sampling highlights **topical** similarities.

Conclusions

- ▶ Argued that the **window-based** contexts of the Paragraph Vector model may have detrimental effect on the learned document embeddings.
- ▶ Proposed a **Context Sampling Framework** that allows for the instantiation of context **policies** of varying complexity.
- ▶ Achieved **significant improvements** over PV on multiple tasks & datasets.