

Text Mining and Language Analytics

Lecture 6

Word embeddings II

Dr Stamos Katsigiannis

2023-24

Embeddings: TF / TF-IDF

- TF/TF-IDF

- Common baseline model
- The standard solution in Information Retrieval!
- Words represented by a simple function of their frequency and the frequency of other words
- Similar words expected to have similar vectors
- High dimensionality ($|V|$ dimensions)
 - Millions of dimensions for large corpora
- Very **sparse** vectors (most elements are zero)

Similar words keep the same company

- **How to increase density and reduce dimensionality?**

- **Solution** → Learn shorter and denser vectors

Embeddings: Sparse vs. Dense vectors

Why dense vectors?

- Dense vectors are typically shorter → Easier to use as features in machine learning
- Dense vectors may generalise better than explicit word frequencies
- Dense vectors may do better at capturing synonymy
- **In practice, they work better!**

Embeddings: Dimensionality reduction

- **Given a Word-Word matrix, how to reduce dimensions?**
 - Filtering
 - Feature selection methods
 - Project data on different space
 - Rotate the axes into a new space where the highest order dimension captures the most variance in the data
 - The next dimension captures the next most variance
 - And so on ...
- Many methods exist, e.g.:
 - **Principal Component Analysis (PCA)**
 - **Singular Value Decomposition (SVD)**
- How to select the number of dimensions?
 - **Empirically decided!**
 - Typically between 50 and 500 dimensions
 - 300 dimensions often used in the literature

Singular Value Decomposition (SVD)

- Consider a Word-Word matrix M
- SVD of M : $M = U\Sigma V^T$

U: Rows correspond to original. Mapped in a different space. Columns ordered by descending variance they account for
 Σ : Diagonal matrix of singular values expressing importance of each dimension

General SVD:

$$\begin{bmatrix} M \\ m \times n \end{bmatrix} = \begin{bmatrix} U \\ m \times m \end{bmatrix} \begin{bmatrix} \Sigma \\ m \times n \end{bmatrix} \begin{bmatrix} V^T \\ n \times n \end{bmatrix}$$

Each row i of U is an embedding of the i^{th} word

SVD on Word-Word matrix:

$$\begin{bmatrix} M \\ |V| \times |V| \end{bmatrix} = \begin{bmatrix} U \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} \Sigma \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} V^T \\ |V| \times |V| \end{bmatrix}$$

Keep top k dimensions

SVD on Word-Word matrix:

$$\begin{bmatrix} M \\ |V| \times |V| \end{bmatrix} = \begin{bmatrix} U \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} \Sigma \\ |V| \times |V| \end{bmatrix} \begin{bmatrix} V^T \\ |V| \times |V| \end{bmatrix}$$

k k k k

Approximation of M

Use U as the embeddings

Usually $k \in [50, 500]$. Also $k = 300$ very common!

Embeddings: word2vec (I)

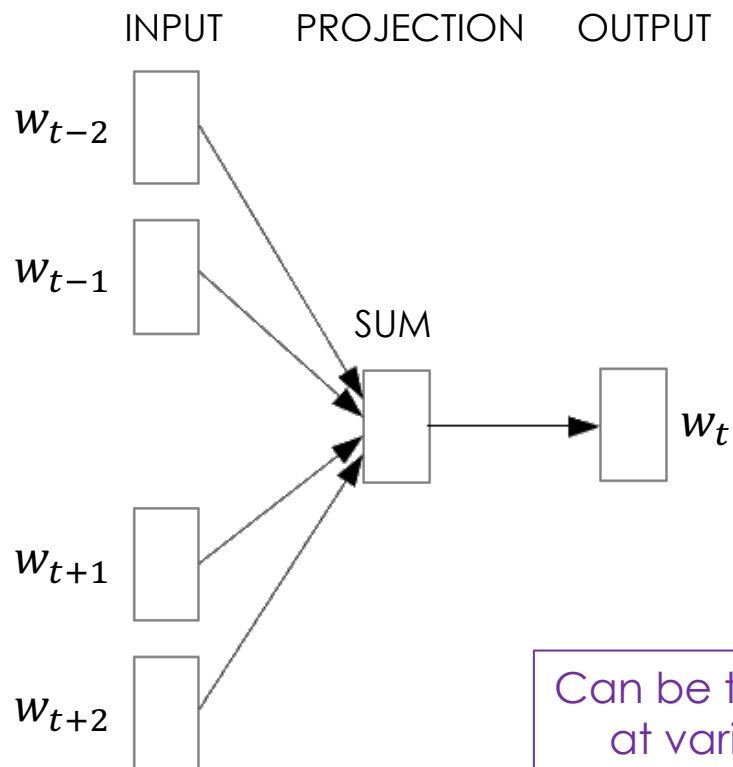
- **Idea** → **Predict than count!**
- Mikolov et al., <https://code.google.com/archive/p/word2vec/>
- Popular embedding method
- Very fast to train
- Simple neural network with one hidden layer
- word2vec is **not a single algorithm**

Embeddings: word2vec (II)

- word2vec **is a software package** for representing words as vectors:
 - Two distinct models
 - Continuous Bag of Words (CBow)
 - Skip-gram
 - Various training methods
 - Negative sampling
 - Hierarchical Softmax
 - Pre-processing pipeline
 - Dynamic Context Windows
 - Subsampling
 - Deleting rare words

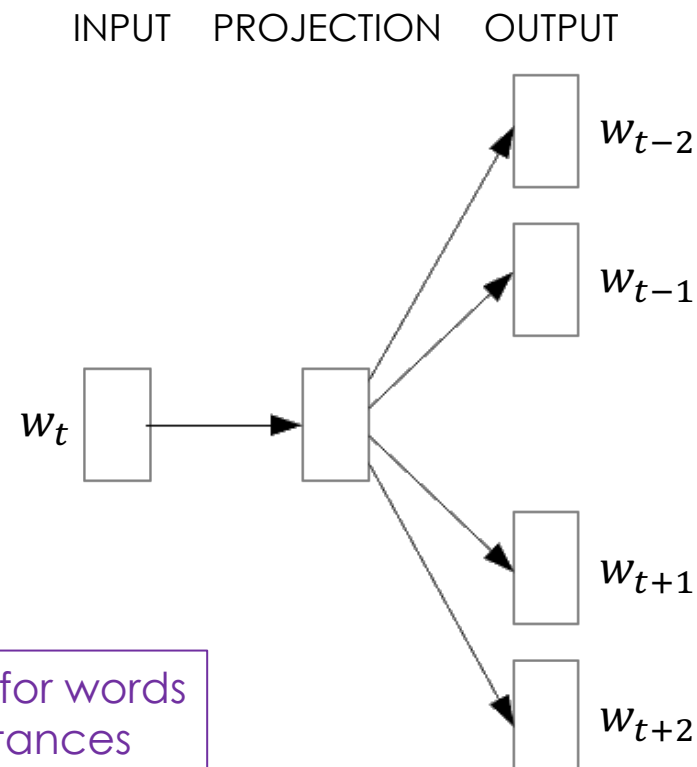
Embeddings: word2vec (III)

CBOW
Predict word
from
surrounding
words



CBoW

Can be trained for words
at various distances

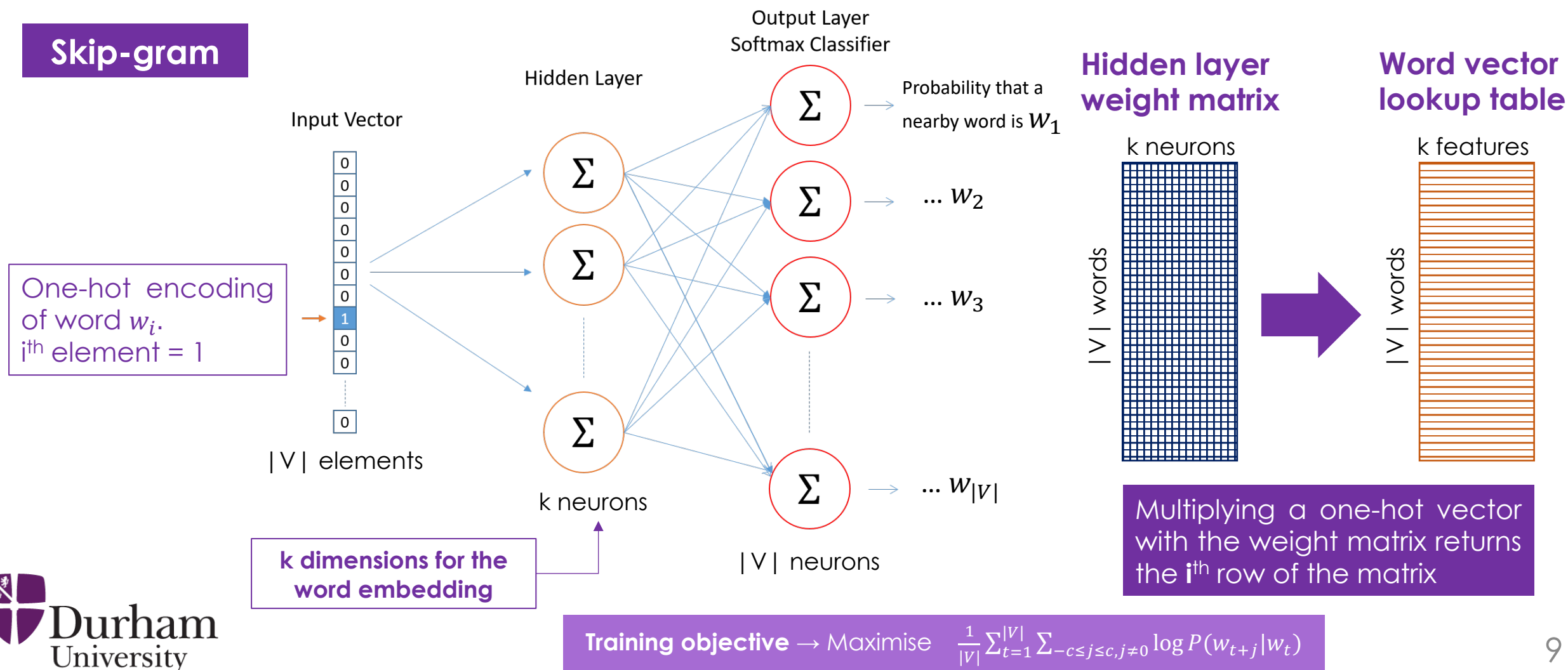


Skip-gram
Use word to
predict
surrounding
words

Skip-gram

Embeddings: word2vec (IV)

Skip-gram



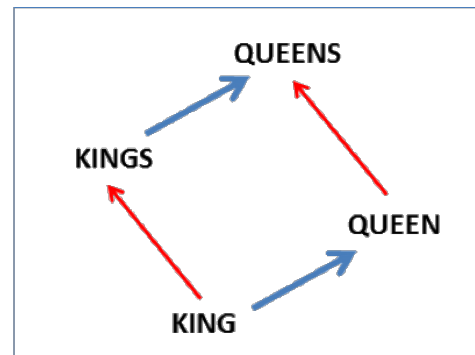
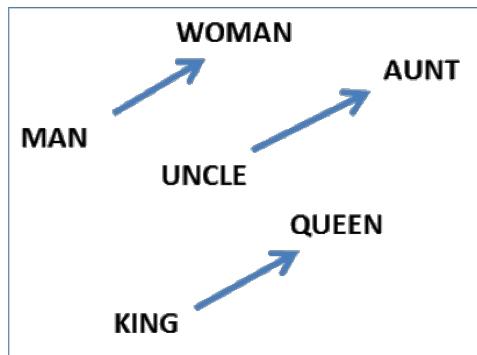
Embeddings: word2vec (V)

Skip-gram training

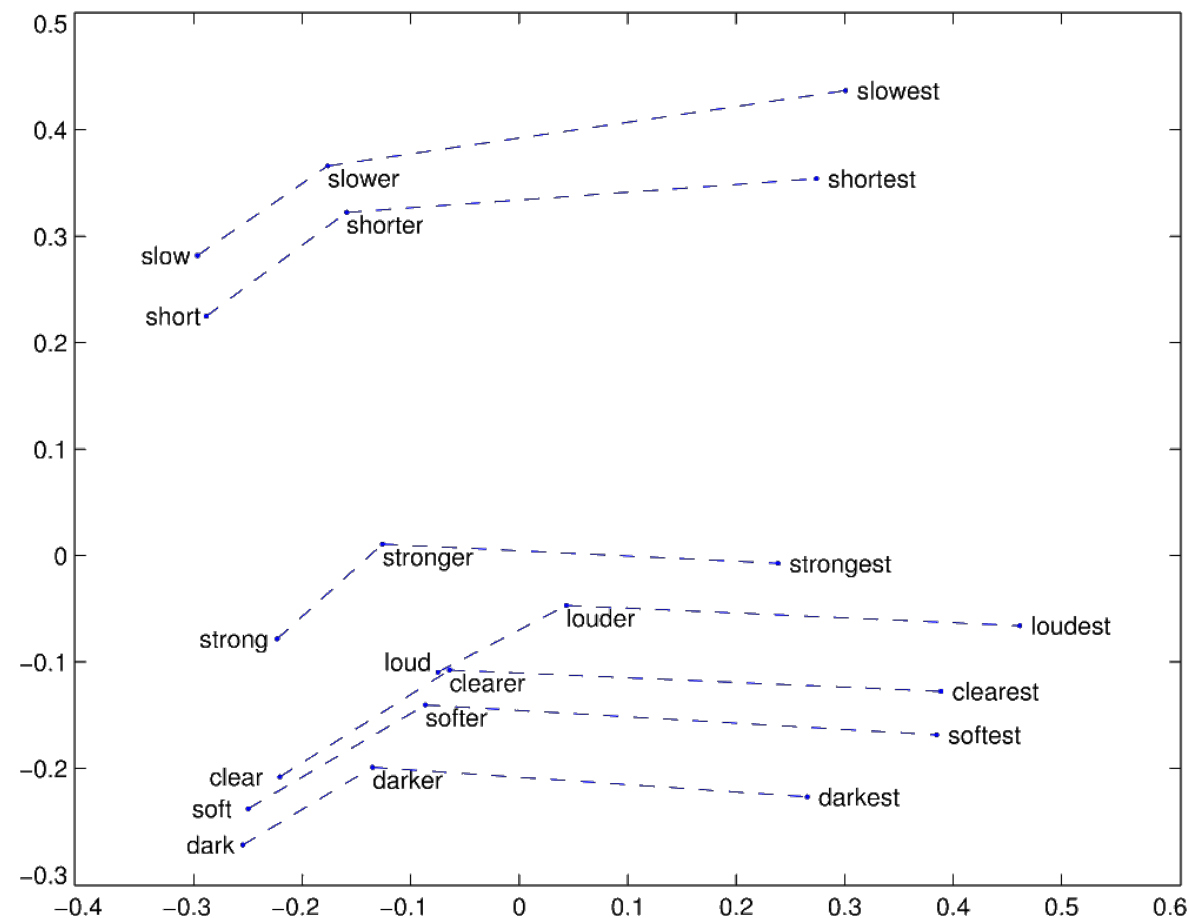
1. Convert words to vectors using one-hot encoding \rightarrow Vector size $1 \times |V|$
 2. Use input layer of $|V|$ neurons to pass word w_t to hidden layer
 3. Hidden layer of k neurons performs dot product between weight matrix $W_{|V| \times k}$ and the input vector $w_t \rightarrow$ The t^{th} row of $W_{|V| \times k}$ will be the output $H_{1 \times k}$
 4. No activation function used $\rightarrow H_{1 \times k}$ passes directly to the output layer
 5. Output layer applies dot product between $H_{1 \times k}$ and $W'_{k \times |V|} \rightarrow$ Output $U_{1 \times |V|}$
 6. Use softmax to compute probability of each word. The word with the highest probability is the result.
 7. If the predicted word for a given context position is wrong, then use backpropagation to modify the weight matrices W and W'
- Execute these steps for each word w_t in vocabulary and each context word
 - Each word w_t will be passed $2 \times c$ times, with c being the size of context
 - Forward propagation will be processed $|V| \times 2 \times c$ times per epoch

Embeddings: Properties (I)

- Similarity depends on context size c
 - e.g., Nearest words to **Hogwarts**
 - $c=\pm 2$: Sunnydale, Evernight
 - $c=\pm 5$: Dumbledore, Malfoy, halfblood
- Embeddings capture relational meaning
 - e.g., $\text{vector}(\text{king}) - \text{vector}(\text{man}) + \text{vector}(\text{woman}) \approx \text{vector}(\text{queen})$
 - e.g., $\text{vector}(\text{Paris}) - \text{vector}(\text{France}) + \text{vector}(\text{Italy}) \approx \text{vector}(\text{Rome})$

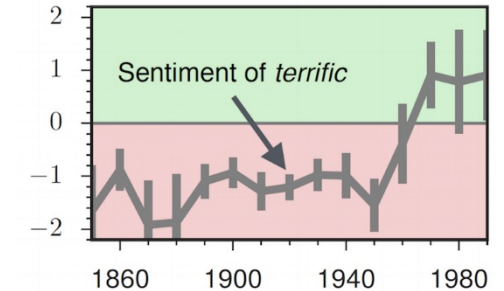


Embeddings: Properties (II)



Embeddings: Properties (III)

- Embeddings can help study word history
 - Train embeddings on old books to study changes in word meaning
 - Remember the N-gram example?
- **Embeddings reflect cultural bias!**
 - $\text{vector}(\text{man}) - \text{vector}(\text{woman}) \approx \text{vector}(\text{computer programmer}) - \text{vector}(\text{homemaker})$
 - Male names associated more with math, female names with arts
 - Words for specific ethnicities associated with stereotypes
 - Embeddings for competence adjectives are biased toward men → Bias slowly decreasing
- Embeddings can be used as historical tool to study bias
 - e.g., The cosine similarity of embeddings for decade X for occupations to male vs female names is correlated with the actual percentage of females and males in that occupation in decade X



- William L. Hamilton, Kevin Clark, Jure Leskovec, Dan Jurafsky, SocialSent: Domain-Specific Sentiment Lexicons for Computational Social Science, <https://nlp.stanford.edu/projects/socialsent/>
- Bolukbasi, Tolga, Kai-Wei Chang, James Y. Zou, Venkatesh Saligrama, and Adam T. Kalai. "Man is to computer programmer as woman is to homemaker? debiasing word embeddings." In Advances in Neural Information Processing Systems, pp. 4349-4357. 2016.
- Garg, Nikhil, Schiebinger, Londa, Jurafsky, Dan, and Zou, James (2018). Word embeddings quantify 100 years of gender and ethnic stereotypes. Proceedings of the National Academy of Sciences, 115(16), E3635-E3644
- Caliskan, Aylin, Joanna J. Brusson and Arvind Narayanan. 2017. Semantics derived automatically from language corpora contain human-like biases. Science 356:6334, 183-186.

Embeddings: Summary

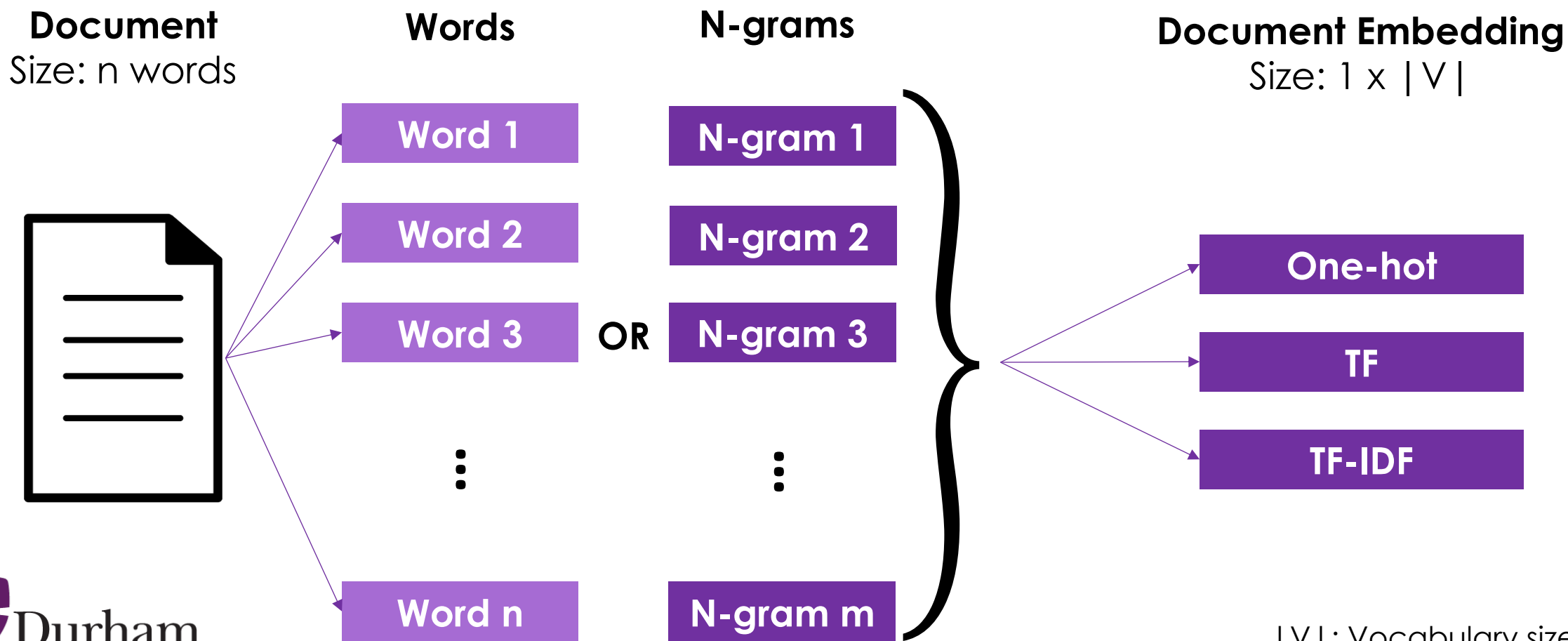
Embeddings → **Vector representations of words that encode meaning**

- More refined representation than a sequence of characters
- Very effective in modelling similarity/analogy → Use cosine distance
- Sparse models (One-hot, TF, TF-IDF, etc.)
- Dense models (word2vec, GloVe, etc.)
- **Work well in practice!**
- Remember that they can encode cultural/dataset bias

Document embeddings

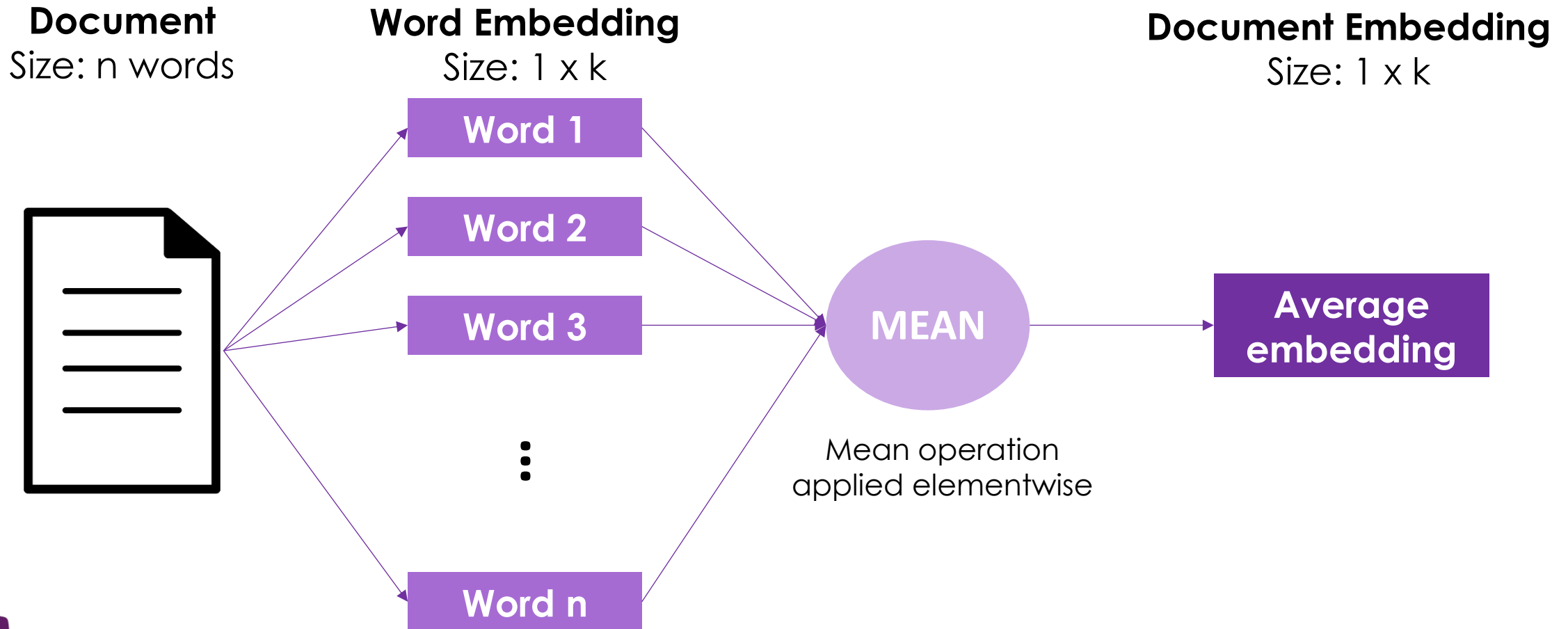
- **What if we wanted to represent whole documents as vectors for machine learning applications?**
 - Topic modelling
 - Sentiment analysis
 - Spam detection
 - Text classification
 - Social media post analysis
 - Product review analysis
 - ...
- **Challenges**
 - Representative information
 - Support high classification accuracy
 - Maintain accuracy across all lengths of documents

Document embeddings: One-hot / TF / TF-IDF / N-gram



|V|: Vocabulary size

Document embeddings: Averaging



Document embeddings: Fixed length (I)

Document
Size: n words



Word Embedding
Size: $1 \times k$

Word 1

Word 2

Word 3

⋮

Word n

n vectors \rightarrow Size: $n \times k$

Word 1 | Word 2 | Word 3 | ... | Word n

3 vectors \rightarrow Size: $3 \times k$

Word 1 | Word 2 | Word 3 | ... | Word n

$n+1$ vectors \rightarrow Size: $(n+1) \times k$

Word 1 | Word 2 | Word 3 | ... | Word n | **Padding**

$n+m$ vectors \rightarrow Size: $(n+m) \times k$

Word 1 | Word 2 | Word 3 | ... | Word n | **Padding** | ... | **Padding**

**Concatenate
word
embeddings!**

Document embedding: Fixed length (II)

- Fixed length is experimental
- **Very short length** → Poor performance (discards valuable information)
- **Very long length** → More computationally expensive
- **Balance between the two needed!**
- Look at the statistics of the training data:
 - **Average** number of words per document in training dataset
 - **Maximum** number of words per document in training dataset
 - **Median** number of words per document in training dataset

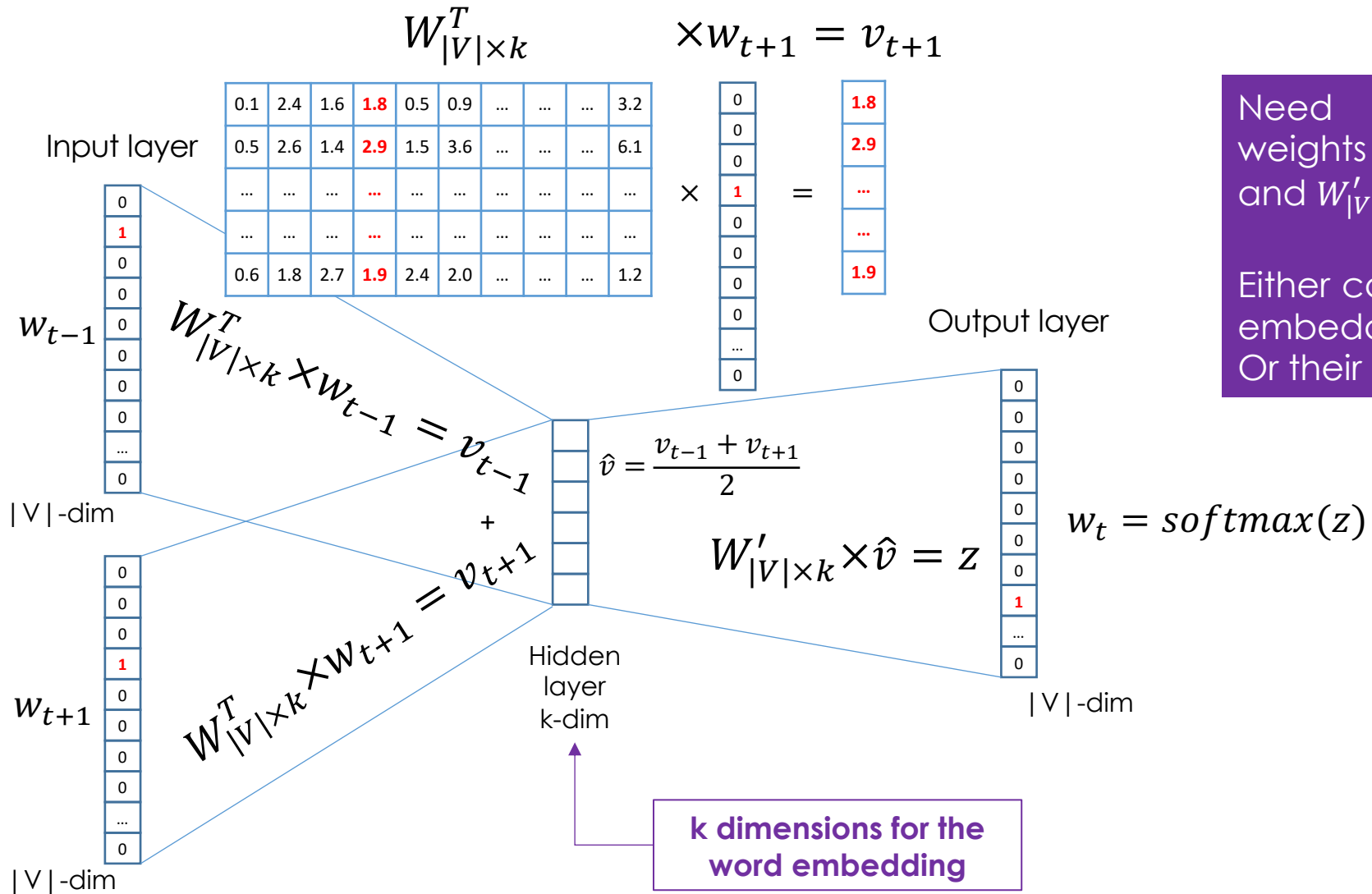
Document embeddings: Summary

- One-hot / TF / TF-IDF –based
 - Sparse vectors
 - High dimensionality
 - Fixed size (typically $|V|$) regardless of document length
- Word embeddings-based
 - Dense vectors
 - High dimensionality
 - Fixed size (k) for averaging approach
 - Fixed size (m) for fixed length approach but may require padding and lead to information loss depending on m and document size

Questions?

Embeddings: word2vec (VI)

CBoW



Need to compute weights matrices $W_{|V| \times k}$ and $W'_{|V| \times k}$

Either can be used as the embedding.
Or their average!