


[WORD EMBEDDINGS] → it's a feature extraction method

Word representation: . 1-hot representation → it doesn't generalize

. factorized = (word embeddings)

→ great for small training set

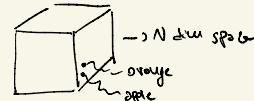
Word Embeddings: TRANSFER LEARNING

1. Learn word emb. from large text corpus (or download pre-trained WE)
2. Train for WE to new task with smaller training set
3. Optional: continue to fine tune the WE with new data

N-features	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

N-dimensional vector to represent the word MAN

↳ it can be represented in a Ndim space:



Properties:

• Reasoning about analogies

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

Annotations for analogy reasoning:

- Man → Woman ≈ King → Queen
- $\frac{e_{\text{Man}} - e_{\text{Woman}}}{e_{\text{Man}} + e_{\text{Woman}}} \approx \frac{e_{\text{King}} - e_{\text{Queen}}}{e_{\text{King}} + e_{\text{Queen}}}$

Usage:

1) Learn an embedding for your problem:

- a) Standalone: word trained to learn an embedding, which is saved and then used as part of another model for your task later

- b) jointly: WE learned as part of a large task-specific model

2) Reuse embeddings (pre-trained embeddings):

- a) static: embedding kept static, used as component of a model
- b) updated: updated during training of your model (fine tuning)

How to learn WE

The distributed representation is learned based on the usage of words. This allows words that are used in similar ways to result in having similar representations, naturally capturing their meaning. This can be contrasted with the crisp and fragile representation in a bag of words model where, unless explicitly managed, different words have different representations, regardless of how they are used.

- Word2Vec: statistical method to learn WE.

CBOW vs SKIP-GRAM Model are 2 learning methods.

Learn about words using context. The context is defined by a window of surrounding words. The window size is a parameter of the model.

- Glove (global vector for word representation):

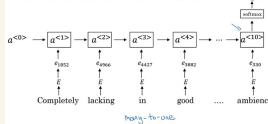
it doesn't use context to define local context.
it builds a word-cooccurrence matrix using statistics across the whole text corpus.

[SENTIMENT CLASSIFICATION]

Problem: small labelled dataset

Solution: word embeddings + RNN

RNN for sentiment classification



[Embedding layer output]

$W = \text{batch size} / \text{training set size}$

$ml = \max \text{len} (\text{sentence max len}) \longrightarrow$

$fs = \text{feature space} (\text{dimension of } W_e)$

