

---

# Text / Word Vectors

# Motivation

- Have shown how to use Neural Networks with structured numerical data.
- Images can be upsampled / downsampled to be a certain size.
- Image values are numbers (greyscale, RGB)
- But how do we work with text?
- Issue 1: How to deal with pieces of text (sequences of words that vary in length)?
- Issue 2: How to convert words into something numerical?

## Issue: Variable length sequences of words

- With images, we forced them into a specific input dimension
- Not obvious how to do this with text
- We will use a new structure of network called a “Recurrent Neural Network” which will be discussed next lecture

# Tokenization

- Need to convert word into something numerical
- First approach: Tokenization
- Treat as a categorical variable with huge number of categories (one hot encoding)
- Deal with some details around casing, punctuation, etc.

“The cat in the hat.”



[ 'the', 'cat', 'in', 'the', 'hat', '<EOS>' ]

# Tokenization

- Use tokens to build a vocabulary
- Vocabulary is a one-to-one mapping from index # to a token
- Usually represented by a list and a dictionary.

index → word

```
[  
  '<EOS>',  
  'the',  
  'cat',  
  'in',  
  'hat',  
  '.',  
]
```

word → index

```
{  
  '<EOS>': 0,  
  'the': 1,  
  'cat': 2,  
  'in': 3,  
  'hat': 4,  
  '.': 5  
}
```

# Issues with Tokenization

- Tokenization loses a lot of information about words:
  - Part of speech
  - Synonymy (distinct words with same or similar meaning)
  - Polysemy (single word with multiple meanings)
  - General context in which word is likely to appear (e.g. “unemployment” and “inflation” ) are both about economics
- Increasing vocabulary size is difficult (would require re-training the model)
- Vector length is huge -> large number of weights
- Yet information in vector is very sparse

# Word Vectors

- Goal: represent a word by an  $m$ -dimensional vector (for medium-sized  $m$ , say,  $m=300$ )
- Have “similar” words be represented by “nearby” vectors in this  $m$ -dimensional space
- Words in a particular domain (economics, science, sports) could be closer to one another than words in other domains.
- Could help with synonymy
  - e.g. “big” and “large” have nearby vectors
- Could help with polysemy
  - “Java” and “Indonesia” could be close in some dimensions
  - “Java” and “Python” are close in other dimensions

# Word Vectors

- Vectors would be shorter length and information-dense, rather than very long and information-sparse
- Would require fewer weights and parameters
- Fortunately, there are existing mappings which can be downloaded and used
- These were trained on big corpora for a long time
- Let's understand how they were developed and trained



# What makes two words similar?

- Idea: similar words occur in similar contexts
- For a given word, look at the words in a “window” around it.
- Consider trying to predict a word given the context.
- This is exactly the CBOW (continuous bag of words) model

“We hold these truths to be **self-evident**, that all men are created equal”



Window size = 3

([‘truths’, ‘to’, ‘be’, ‘that’, ‘all’, ‘men’], ‘self-evident’)



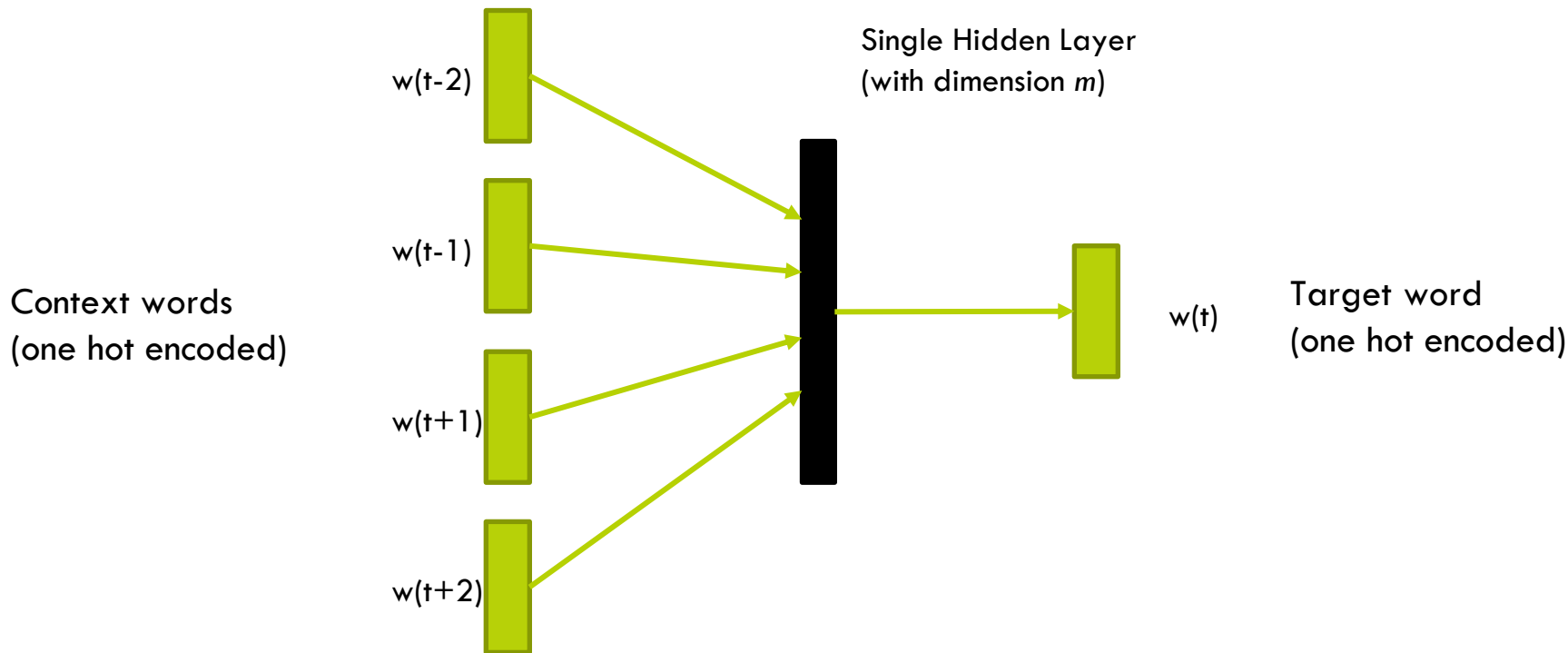
context



target word

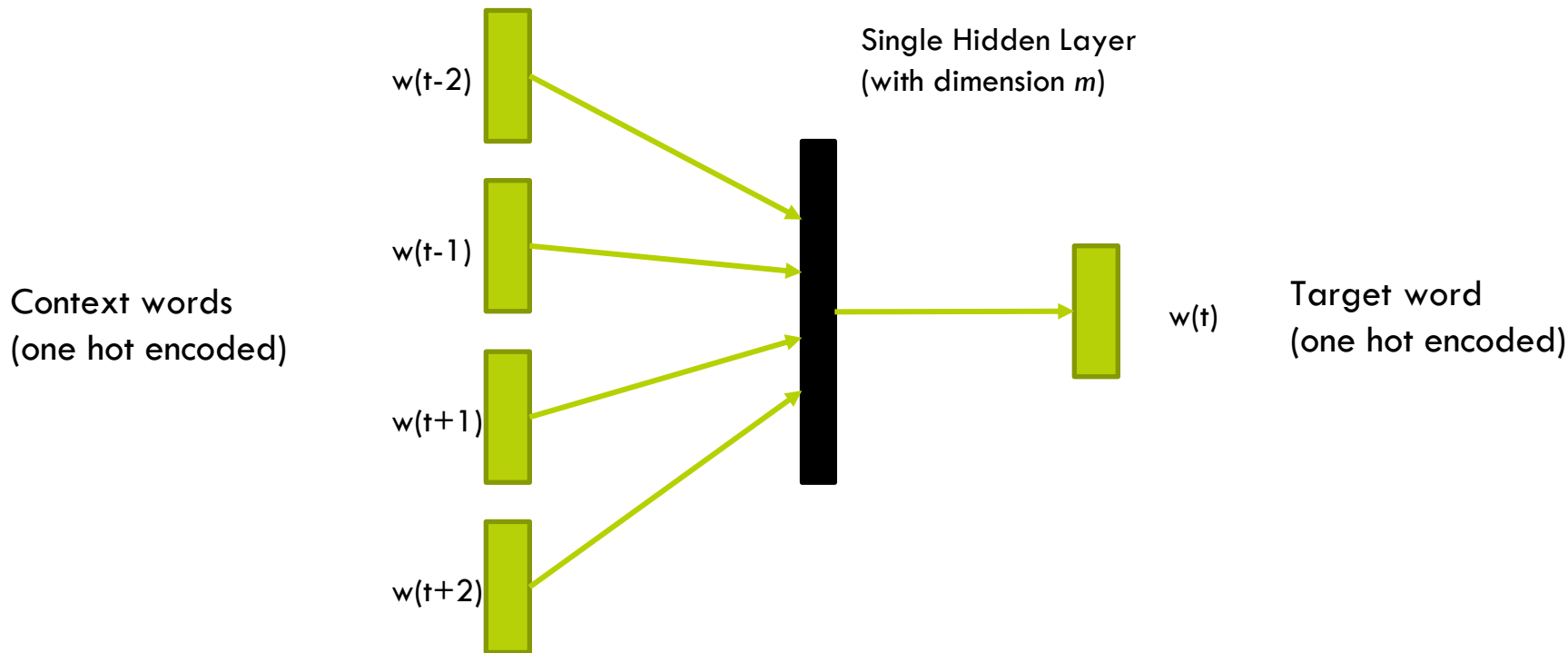
# CBOW Model

Train a neural network on a large corpus of data



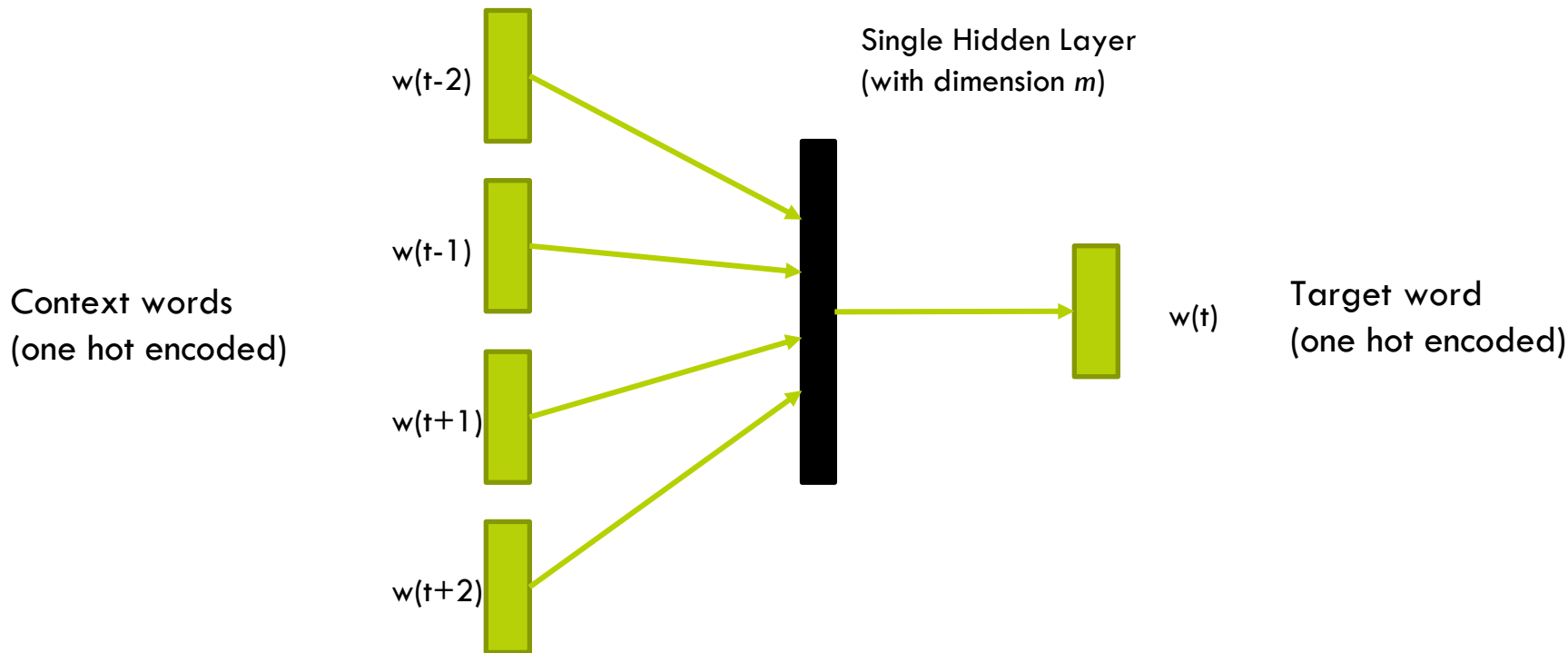
# CBOW Model

Once the network is trained, weights  $\rightarrow$  word vectors



# CBOW Model

Once the network is trained, weights  $\rightarrow$  word vectors



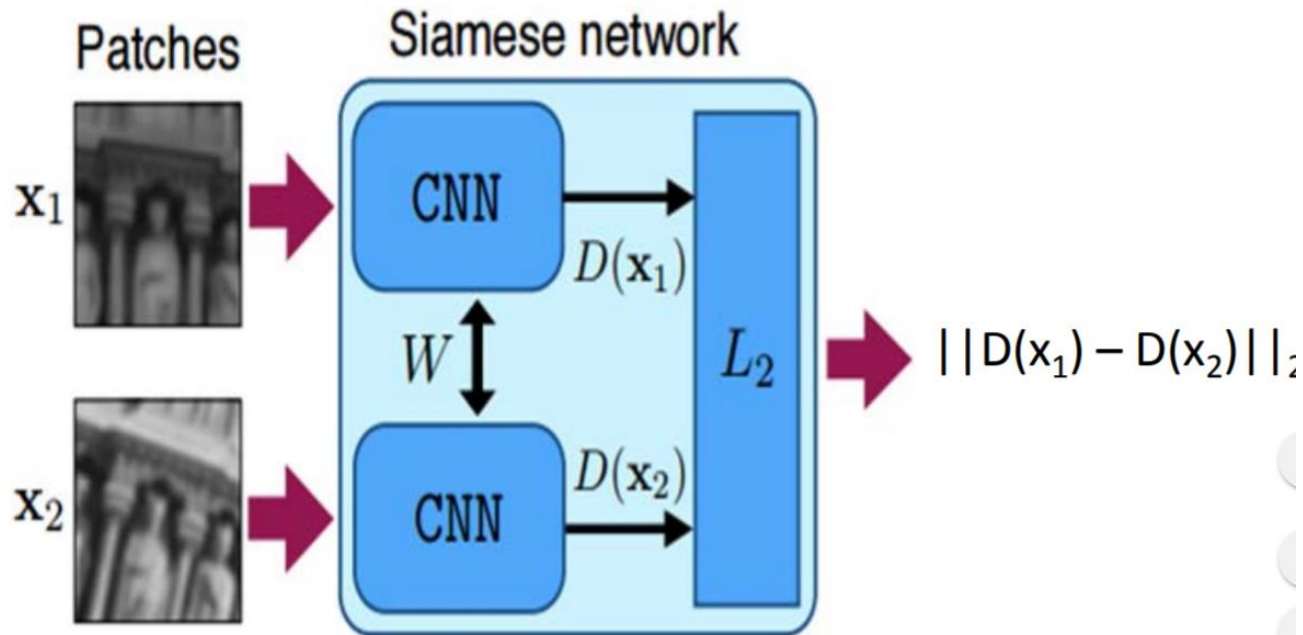
# Other Methods

SkipGram

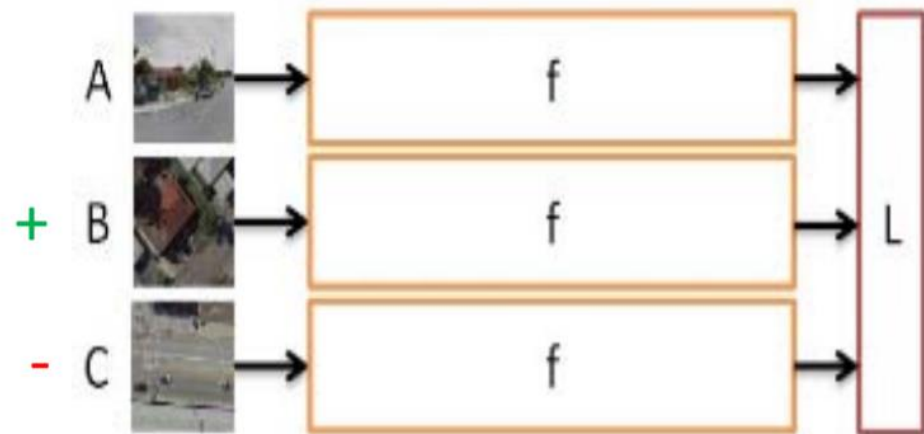
Glove

# Face Verification

Viola Jpnes algo for face detection



# TRIPLET NETWORK



$$\mathbf{D}(f(A), f(B)) < \mathbf{D}(f(A), f(C))$$

# GloVe

- GloVe is publicly available
- Developed at Stanford:  
<https://nlp.stanford.edu/projects/glove/>
- Trained on huge corpora



# GloVe

- Global Vectors for Word Representation (GloVe)
- Use co-occurrence matrix with neighboring words to determine similarity

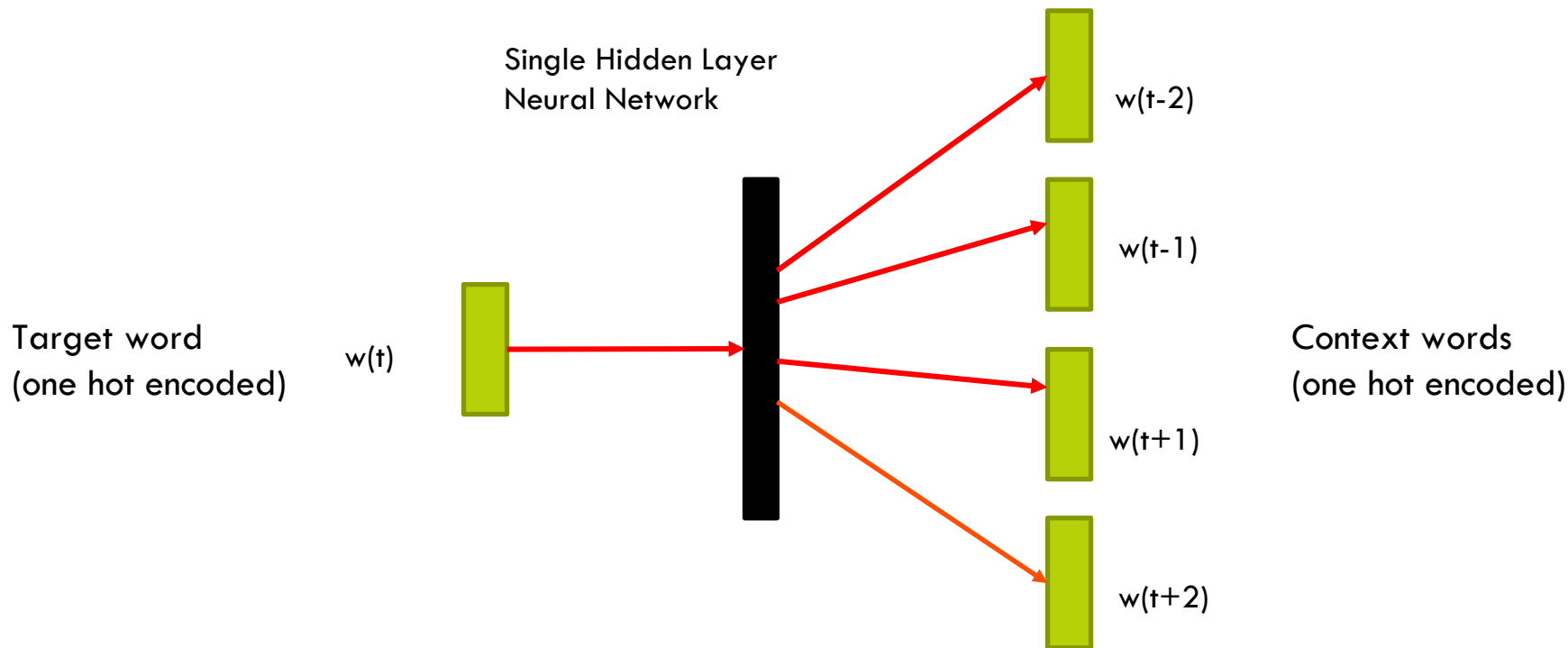
$$J = \frac{1}{2} \sum_{i,j=1}^W f(P_{ij}) (u_i^T v_j - \log(P_{ij}))^2$$

$f \rightarrow$  frequency of a word, with a maximum cap

$P_{ij} \rightarrow$  probability words  $i$  and  $j$  occur together

# Skip-gram Model

Same idea, except we predict the context from the target.



# Word2Vec

- *Distributed Representations of Words and Phrases and Their Compositionality* – Mikolov et al.
- Uses a Skip-gram model to train on a large corpus
- Lots of details to make it work better
  - Aggregation of multi-word phrases (e.g. Boston Globe)
  - Subsampling (i.e. oversample less common words)
  - Negative Sampling (give network examples of wrong words)