

Word2Vec

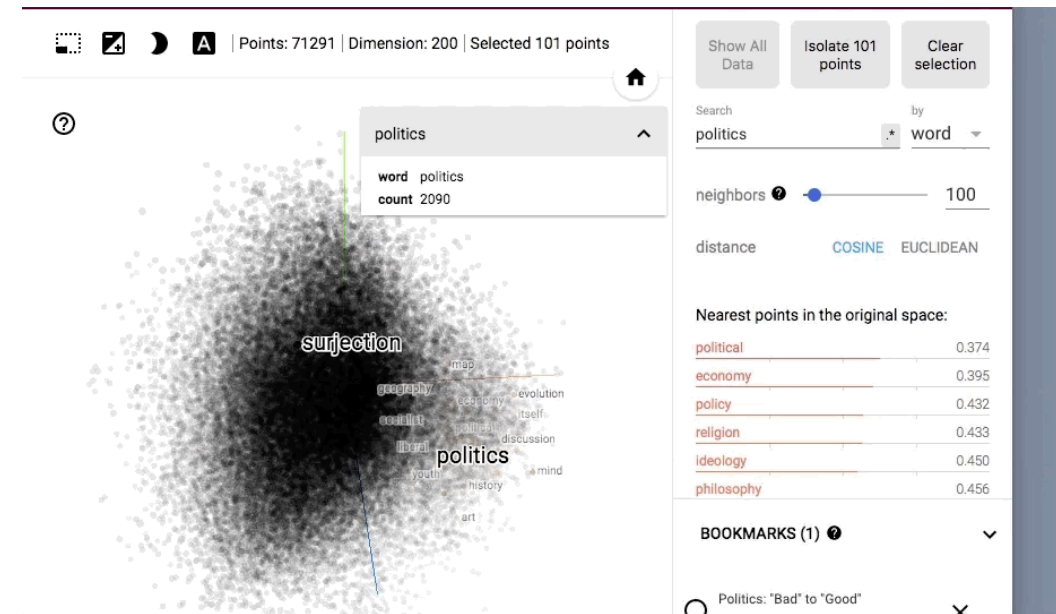
Outline

- Motivation
- Word2Vec from Scratch
- The 20 Newsgroup Text Dataset

Motivation

Word2Vec

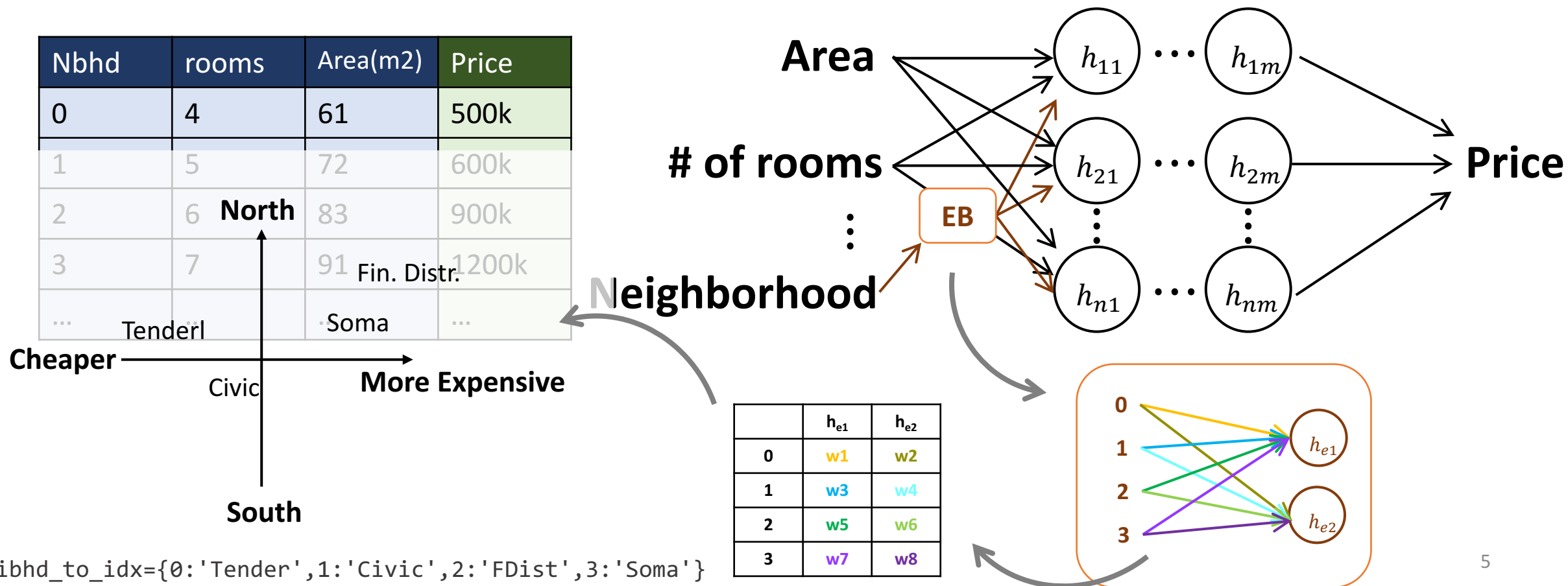
- **Input:** Wikipedia corpus
- **Dimension:** 200
- **Output (CBoW):** The model predicts the current word from a window of surrounding context words
- **Output (CSG):** The model uses the current word to predict the surrounding window of context words



Examples from: <http://projector.tensorflow.org/>

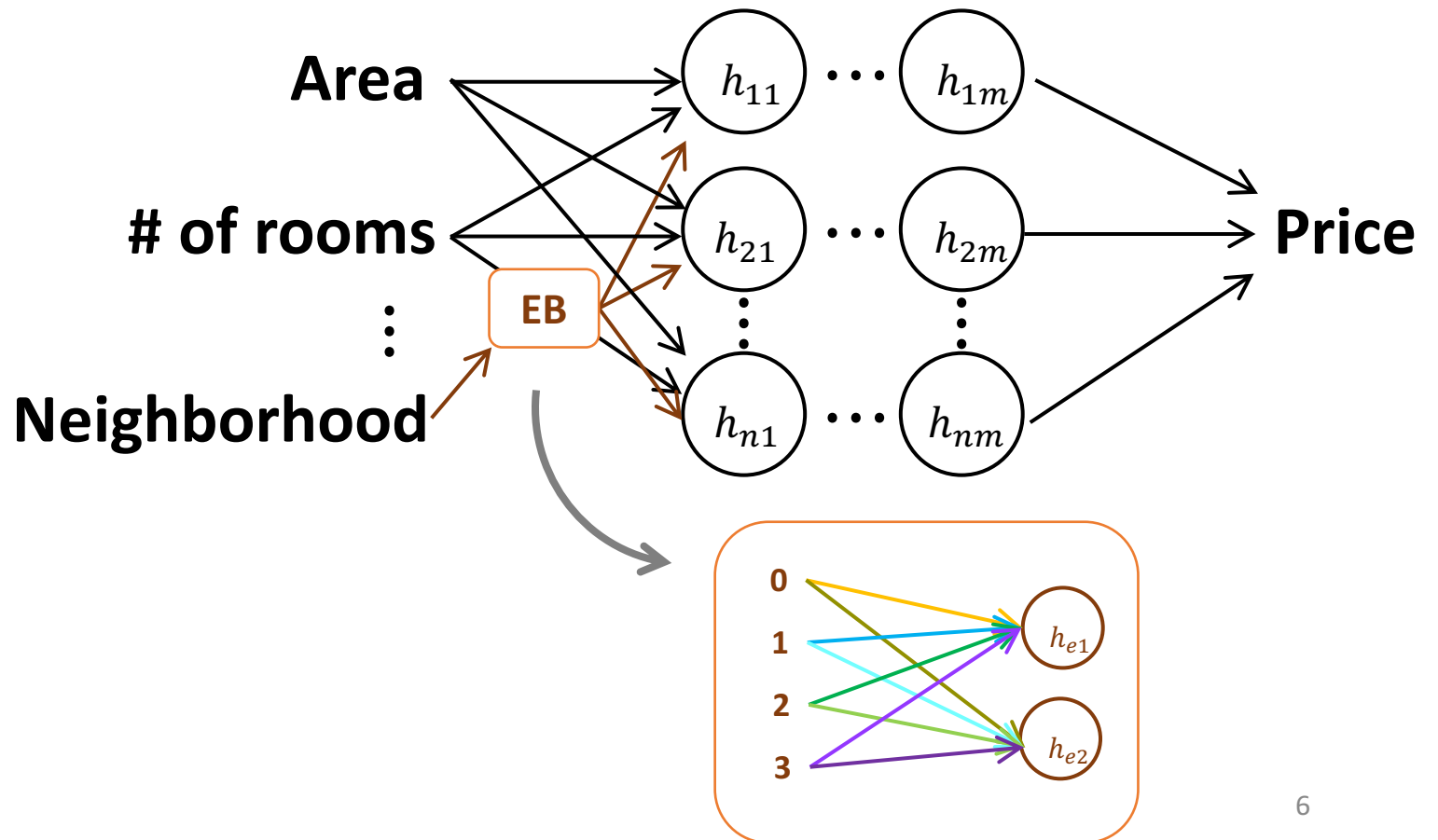
Embeddings

- Recap from last Class



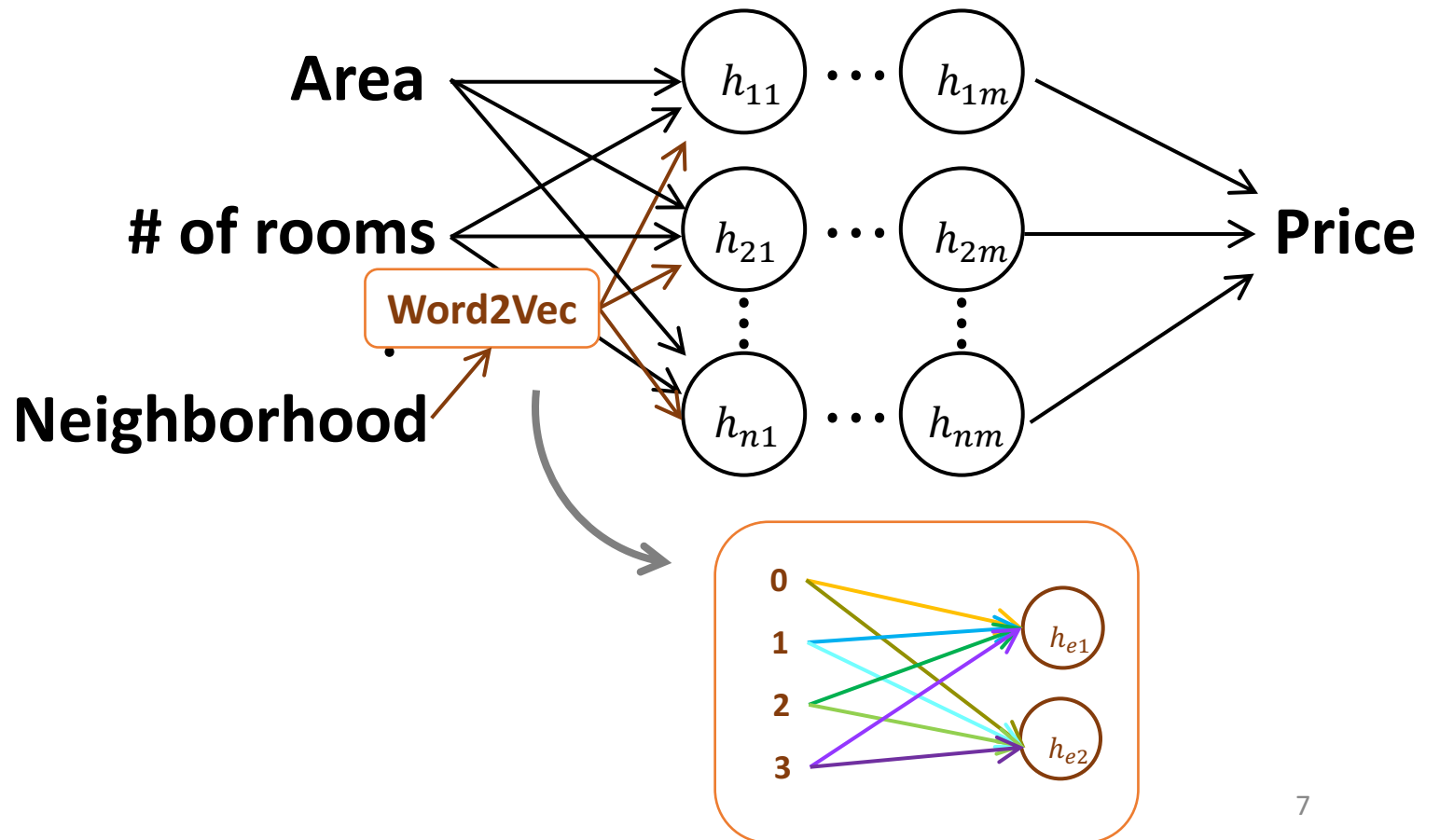
Embeddings

- Recap from last Class



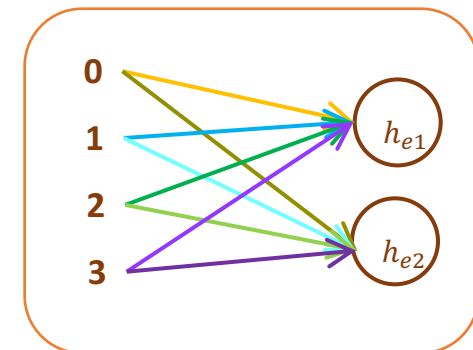
Embeddings

- Recap from last Class



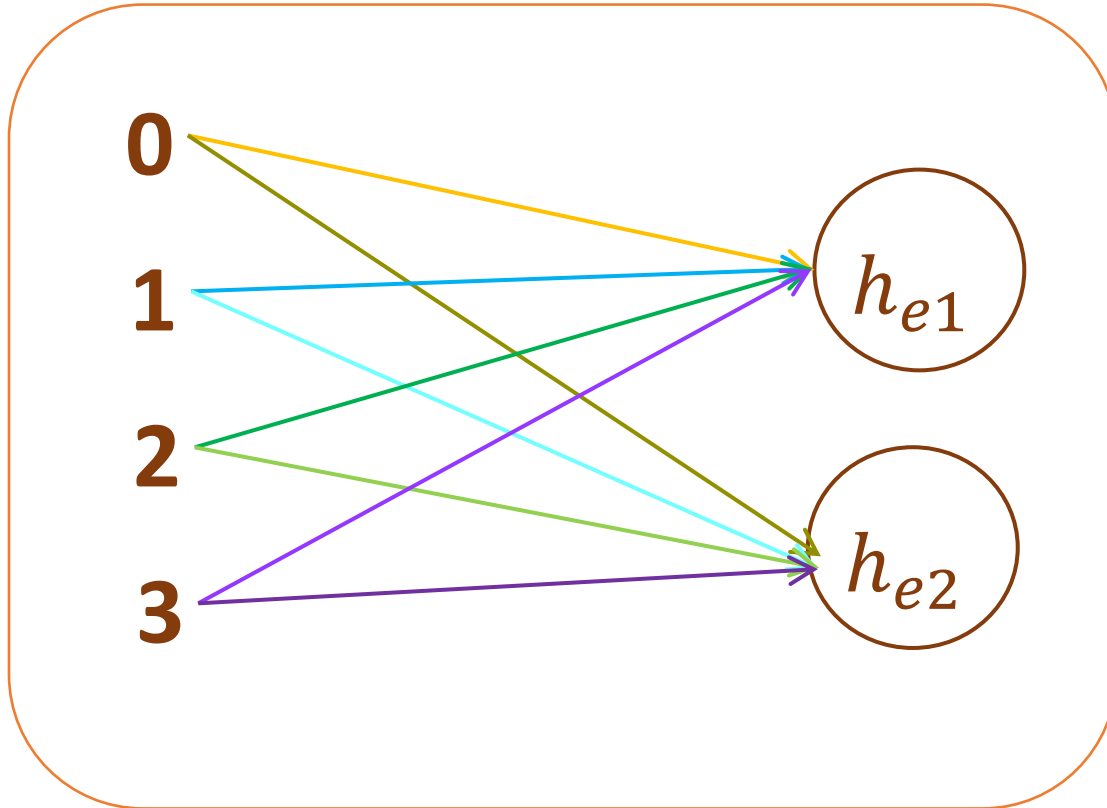
Embeddings

- Recap from last Class



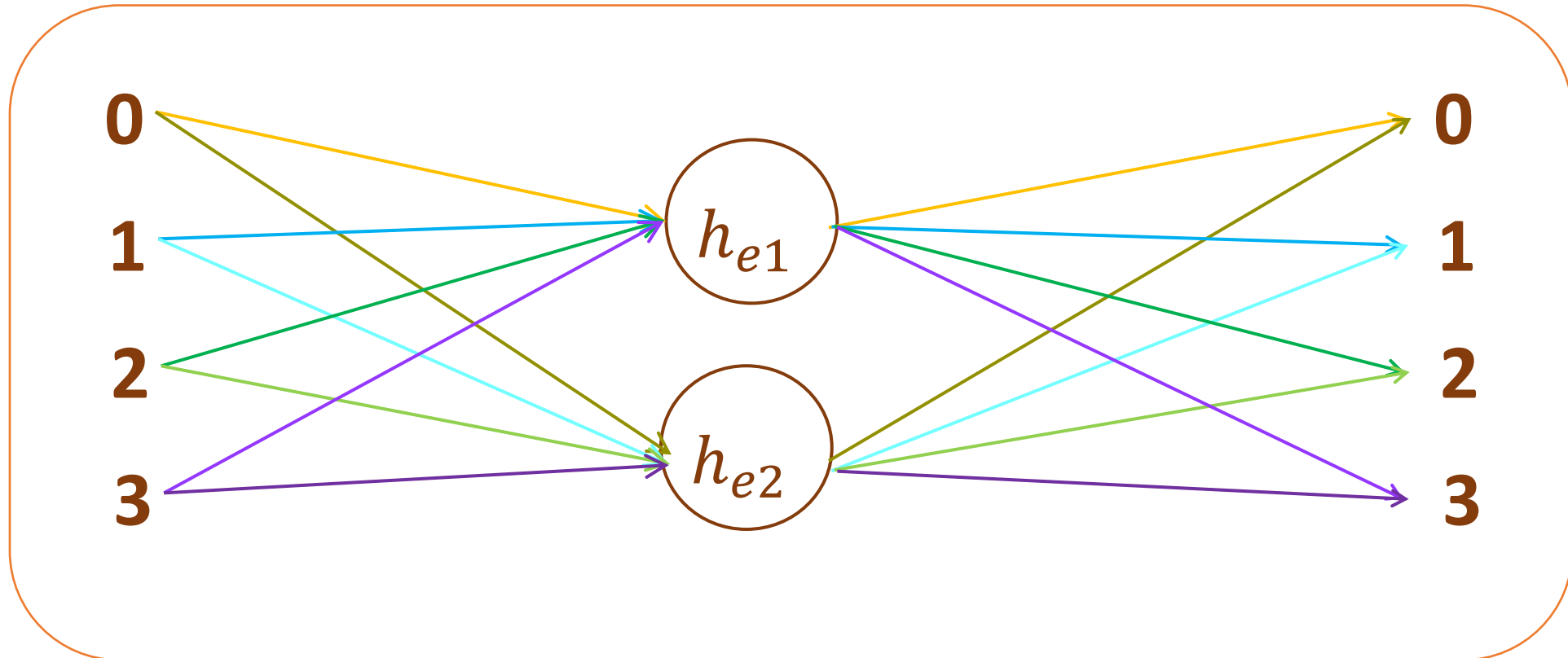
Embeddings

- Recap from last Class



Embeddings

- Recap from last Class

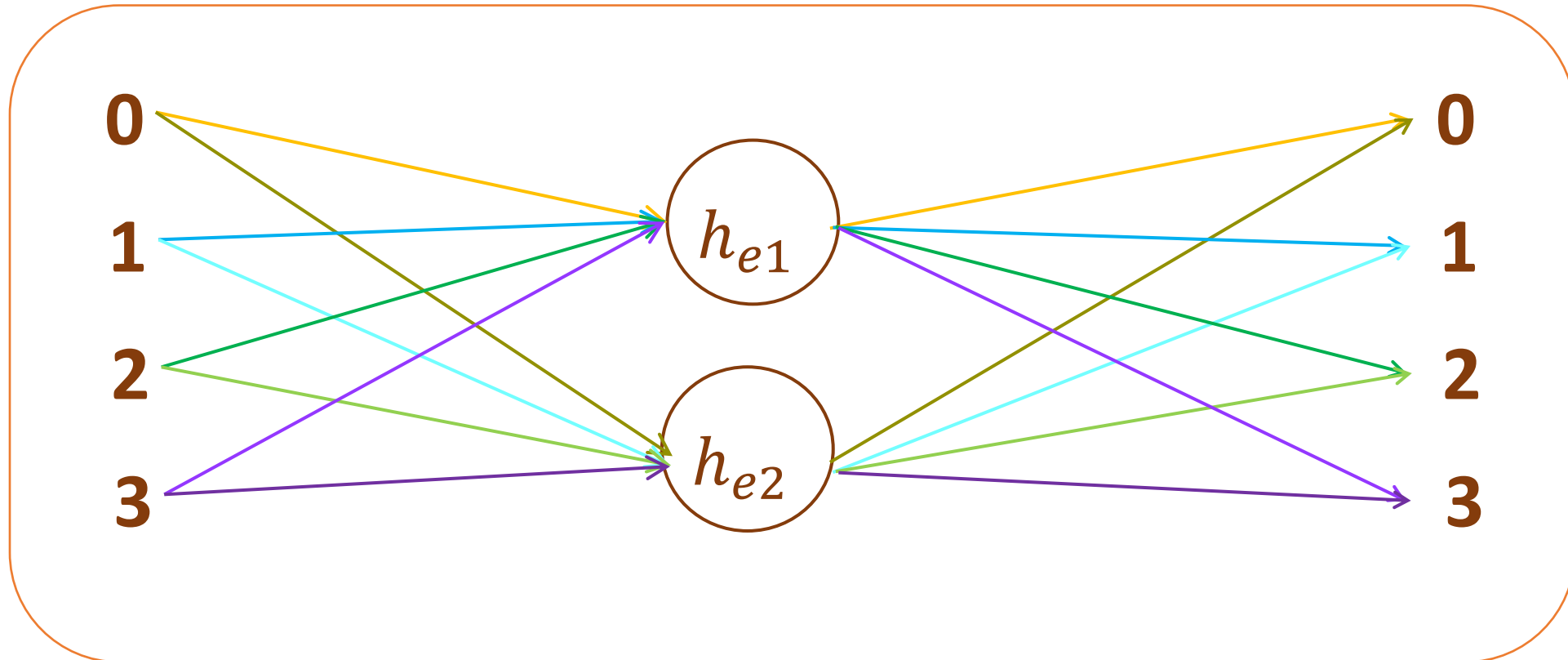


Embeddings

- Recap from last Class

I am the king

0 1 2 3

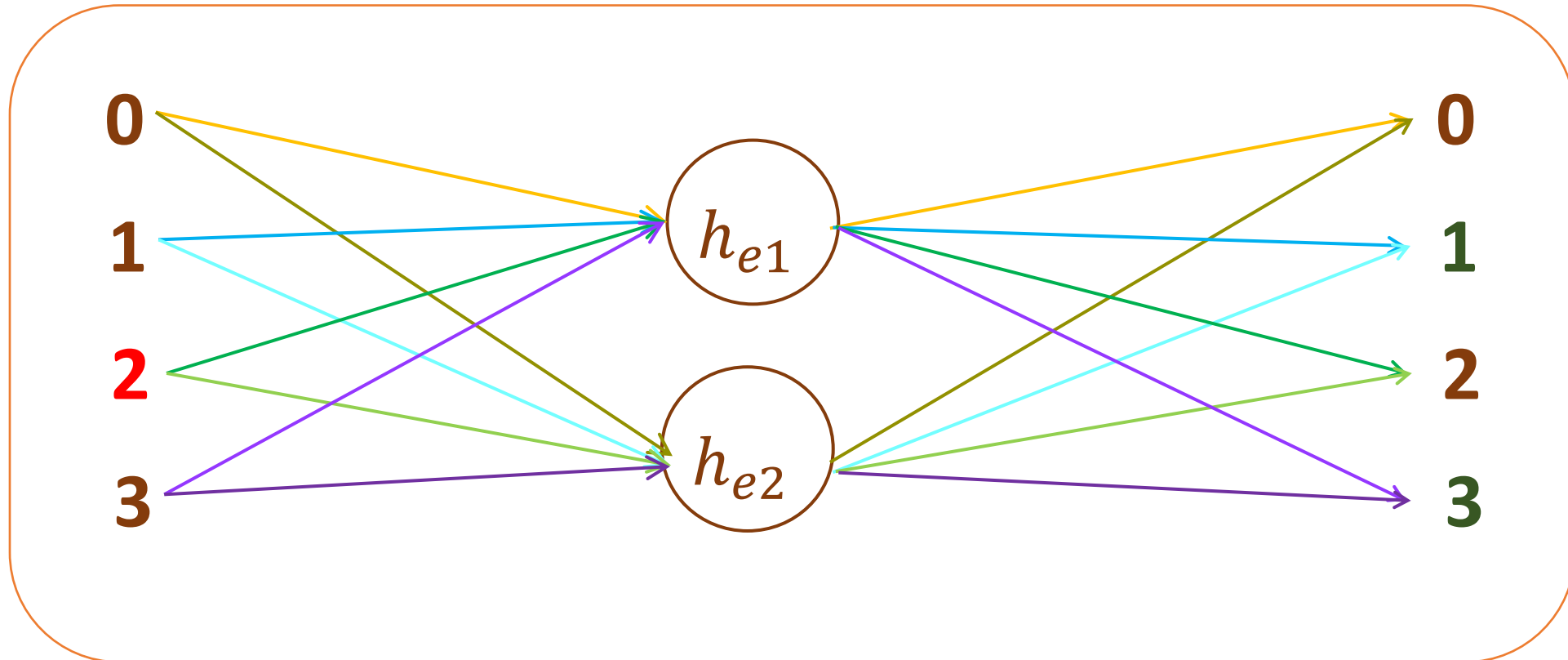


Embeddings

- Recap from last Class

I am **the** king

0 1 2 3



Embeddings

- How to build the dataset?
 - <https://www.youtube.com/watch?v=64qSgA66P-8>

Word2Vec data generation (skipgram)
(window size = 1)

“king brave man”
“queen beautiful woman”

word	neighbor
king	brave
brave	king
brave	man
man	brave
queen	beautiful
beautiful	queen
beautiful	woman
woman	beautiful

Embeddings

- How to build the dataset?
 - <https://www.youtube.com/watch?v=64qSgA66P-8>

Word2Vec data generation (skip gram)
(window size = 2)

"king brave man"
"queen beautiful woman"
✱

word	neighbor
king	brave
king	man
brave	king
brave	man
man	king
man	brave
queen	beautiful
queen	woman
beautiful	queen
beautiful	woman
woman	queen
woman	beautiful

Word2Vec from Scratch