# Latent Semantic Indexing and Word Vector Representations
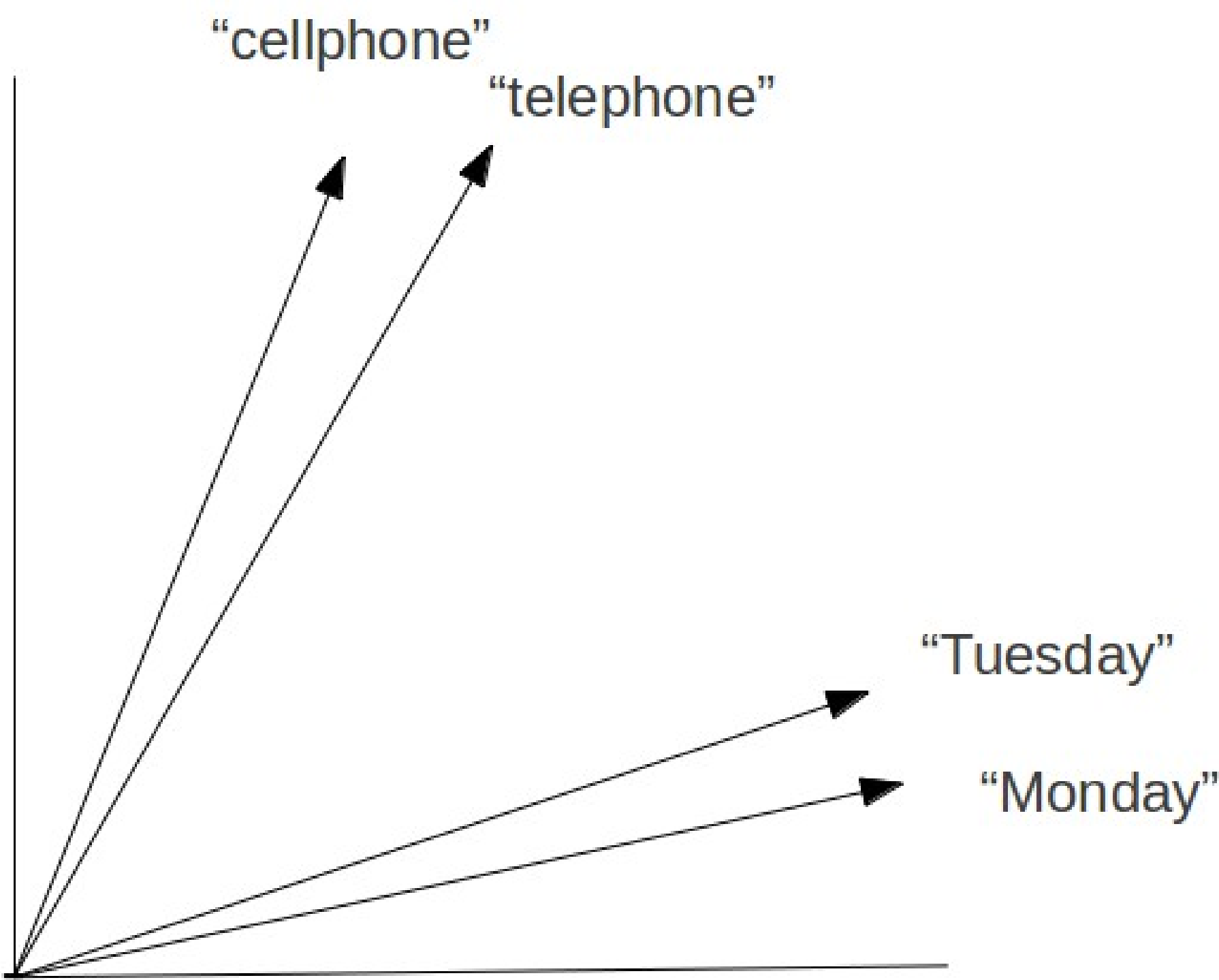
Trevor Barron

Adviser: Matthew Whitehead – Colorado College

## Introduction

- In a probabilistic *N*-gram natural language model each word is represented as a single element. A word's representation does not distinguish it from any other word and all words are equally related. "Monday" and "telephone" are just as related as "Monday" and "Tuesday" in this model.
- Representing words using semantic feature vectors in an *N* dimensional space permits comparison of words by the distance between their respective vectors. The distance between the vectors representing "Monday" and "Tuesday" is much less than the distance between the vectors for "Monday" and "telephone."
- One can find a vector representation for a word by latent semantic indexing (LSI) on a large text corpus.
- This research explores using neural networks to predict the abstract *meaning* of future words during text generation using vector representations for words.

## Objectives

Perform LSI on a large text corpus to generate a vector representation for words in *N* dimensional space. Use these vectors as input and output to train a neural network to understand the mapping between and abstract input and an abstract output.



## Methods

- LSI term document matrix X
  - X(*i, j*) = Count(term *i* in document *j*)
  - Term frequency – inverse document frequency preconditioning of *X* resulted in significant improvements.
- Use Singular Value Decomposition (SVD) to find an approximate vector representation

$$
(\mathbf{t}_i^T) \rightarrow
\begin{bmatrix}
x_{1,1} & \cdots & x_{1,n} \\
\vdots & \ddots & \vdots \\
x_{m,1} & \cdots & x_{m,n}
\end{bmatrix}
= (\mathbf{t}_i^T) \rightarrow
\begin{bmatrix}
\begin{bmatrix} \\ \mathbf{u}_1 \\ \\ \end{bmatrix} \cdots \begin{bmatrix} \\ \mathbf{u}_l \\ \\ \end{bmatrix}
\end{bmatrix}
\cdot
\begin{bmatrix}
\sigma_1 & \cdots & 0 \\
\vdots & \ddots & \vdots \\
0 & \cdots & \sigma_l
\end{bmatrix}
\cdot
\begin{bmatrix}
\begin{bmatrix} & \mathbf{v}_1 & \end{bmatrix} \\
\vdots \\
\begin{bmatrix} & \mathbf{v}_l & \end{bmatrix}
\end{bmatrix}
$$

- Find the k-rank reduction by finding the product of the *k* largest values in Σ and the associated vectors in *U*.
  - The vector for term *ti* = *U*Σ[*i*], the *i*-th row of *U*Σ.

- Use the vectors of dataset trigrams as input and output to train a neural network
  - Example: "the morning sun"
    - the = [-109.2, 5.9, 4.4, 3.8, 1.0]; morning = [-10.1, -8.6, -0.6, 1.6, -0.9]; sun = [-6.9, -4.7, 2.9, 1.6, 0.6]
    - Input = [-109.2, 5.9, 4.4, 3.8, 1.0, -10.1, -8.6, -0.6, 1.6, -0.9]
    - Output = [-6.9, -4.7, 2.9, 1.6, 0.6]

## Results

- Using the Brown corpus (500 documents and 46758 unique terms after preprocessing) to generate vector representations.
  - Sample KNN:



- This works very well for some cases as seen above but very poorly for others. A larger dataset may resolve some of these problems.



- Neural network learning was difficult even for carefully chosen datasets.

## Conclusion

- It takes a lot of effort to find a good vector representation
- SVD is quite sensitive to the condition of the matrix. One very common word like "the" can throw off the entire set of vectors
- In almost every case, more data equals a better result.
- A data formatting problem may be causing the neural network training difficulties. Moving forward one might ensure that the input & output are produced as expected and perhaps try alternative methods of normalization.

## References

Scott Deerwester, Susan T. Dumais, George W. Furnas, Thomas K. Landauer, Richard Harshman (1990). "Indexing by Latent Semantic Analysis" (PDF). Journal of the American Society for Information Science 41 (6): 391–407.

W. N. Francis, H. Kucera, A Standard Corpus of Present-Day Edited American English, for use with Digital Computers (Brown Corpus). 1964.