

Half Title



Title Page

LOC Page

*To my dog
and my cat.*



Contents



Foreword

I am delighted to introduce the first book on Multimedia Data Mining. When I came to know about this book project undertaken by two of the most active young researchers in the field, I was pleased that this book is coming in early stage of a field that will need it more than most fields do. In most emerging research fields, a book can play a significant role in bringing some maturity to the field. Research fields advance through research papers. In research papers, however, only a limited perspective could be provided about the field, its application potential, and the techniques required and already developed in the field. A book gives such a chance. I liked the idea that there will be a book that will try to unify the field by bringing in disparate topics already available in several papers that are not easy to find and understand. I was supportive of this book project even before I had seen any material on it. The project was a brilliant and a bold idea by two active researchers. Now that I have it on my screen, it appears to be even a better idea.

Multimedia started gaining recognition in 1990s as a field. Processing, storage, communication, and capture and display technologies had advanced enough that researchers and technologists started building approaches to combine information in multiple types of signals such as audio, images, video, and text. Multimedia computing and communication techniques recognize correlated information in multiple sources as well as insufficiency of information in any individual source. By properly selecting sources to provide complementary information, such systems aspire, much like human perception system, to create a holistic picture of a situation using only partial information from separate sources.

Data mining is a direct outgrowth of progress in data storage and processing speeds. When it became possible to store large volume of data and run different statistical computations to explore all possible and even unlikely correlations among data, the field of data mining was born. Data mining allowed people to hypothesize relationships among data entities and explore support for those. This field has been put to applications in many diverse domains and keeps getting more applications. In fact many new fields are direct outgrowth of data mining and it is likely to become a powerful computational tool.



Preface

Approximately 17 million people in the USA (6% of the population) and 140 million people worldwide (this number is expected to rise to almost 300 million by the year 2025) suffer from *diabetes mellitus*. Currently, there are a few dozens of commercialised devices for detecting blood glucose levels [1]. However, most of them are invasive. The development of a noninvasive method would considerably improve the quality of life for diabetic patients, facilitate their compliance for glucose monitoring, and reduce complications and mortality associated with this disease. Noninvasive and continuous monitoring of glucose concentration in blood and tissues is one of the most challenging and exciting applications of optics in medicine. The major difficulty in development and clinical application of optical noninvasive blood glucose sensors is associated with very low signal produced by glucose molecules. This results in low sensitivity and specificity of glucose monitoring by optical methods and needs a lot of efforts to overcome this difficulty.

A wide range of optical technologies have been designed in attempts to develop robust noninvasive methods for glucose sensing. The methods include infrared absorption, near-infrared scattering, Raman, fluorescent, and thermal gradient spectroscopies, as well as polarimetric, polarization heterodyning, photonic crystal, optoacoustic, optothermal, and optical coherence tomography (OCT) techniques [1-31].

For example, the polarimetric quantification of glucose is based on the phenomenon of optical rotatory dispersion, whereby a chiral molecule in an aqueous solution rotates the plane of linearly polarized light passing through the solution. The angle of rotation depends linearly on the concentration of the chiral species, the pathlength through the sample, and the molecule specific rotation. However, polarization sensitive optical technique makes it difficult to measure *in vivo* glucose concentration in blood through the skin because of the strong light scattering which causes light depolarization. For this reason, the anterior chamber of the eye has been suggested as a sight well suited for polarimetric measurements, since scattering in the eye is generally very low compared to that in other tissues, and a high correlation exists between the glucose in the blood and in the aqueous humor. The high accuracy of anterior eye chamber measurements is also due to the low concentration of optically active aqueous proteins within the aqueous humor.

On the other hand, the concept of noninvasive blood glucose sensing using the scattering properties of blood and tissues as an alternative to spectral absorption and polarization methods for monitoring of physiological glucose

concentrations in diabetic patients has been under intensive discussion for the last decade. Many of the considered effects, such as changing of the size, refractive index, packing, and aggregation of RBC under glucose variation, are important for glucose monitoring in diabetic patients. Indeed, at physiological concentrations of glucose, ranging from 40 to 400 mg/dl, the role of some of the effects may be modified, and some other effects, such as glucose penetration inside the RBC and the followed hemoglobin glycation, may be important [30-32].

Noninvasive determination of glucose was attempted using light scattering of skin tissue components measured by a spatially-resolved diffuse reflectance or NIR frequency-domain reflectance techniques. Both approaches are based on change in glucose concentration, which affects the refractive index mismatch between the interstitial fluid and tissue fibers, and hence reduces scattering coefficient. A glucose clamp experiment showed that reduced scattering coefficient measured in the visible range qualitatively tracked changes in blood glucose concentration for the volunteer with diabetes studied.

Contributors

Michael Aftosmis

NASA Ames Research Center
Moffett Field, California

Pratul K. Agarwal

Oak Ridge National Laboratory
Oak Ridge, Tennessee

Sadaf R. Alam

Oak Ridge National Laboratory
Oak Ridge, Tennessee

Gabrielle Allen

Louisiana State University
Baton Rouge, Louisiana

Martin Sandve Alnæs

Simula Research Laboratory and
University of Oslo, Norway
Norway

Steven F. Ashby

Lawrence Livermore National
Laboratory
Livermore, California

David A. Bader

Georgia Institute of Technology
Atlanta, Georgia

Benjamin Bergen

Los Alamos National Laboratory
Los Alamos, New Mexico

Jonathan W. Berry

Sandia National Laboratories
Albuquerque, New Mexico

Martin Berzins

University of Utah
Salt Lake City, Utah

Abhinav Bhatele

University of Illinois
Urbana-Champaign, Illinois

Christian Bischof

RWTH Aachen University
Germany

Rupak Biswas

NASA Ames Research Center
Moffett Field, California

Eric Böhm

University of Illinois
Urbana-Champaign, Illinois

James Bordner

University of California, San Diego
San Diego, California

George Bosilca

University of Tennessee
Knoxville, Tennessee

Greg L. Bryan

Columbia University
New York, New York

Marian Bubak

AGH University of Science and
Technology
Kraków, Poland

Andrew Canning

Lawrence Berkeley National
Laboratory
Berkeley, California

Jonathan Carter

Lawrence Berkeley National
Laboratory
Berkeley, California

Zizhong Chen

Jacksonville State University
Jacksonville, Alabama

Joseph R. Crobak

Rutgers, The State University of
New Jersey
Piscataway, New Jersey

Roxana E. Diaconescu

Yahoo! Inc.
Burbank, California

Peter Diener

Louisiana State University
Baton Rouge, Louisiana

Jack J. Dongarra

University of Tennessee, Knoxville,
Oak Ridge National Laboratory,
and
University of Manchester

John B. Drake

Oak Ridge National Laboratory
Oak Ridge, Tennessee

Kelvin K. Droegemeier

University of Oklahoma
Norman, Oklahoma

Stéphane Ethier

Princeton University
Princeton, New Jersey

Christoph Freundl

Friedrich–Alexander–Universität
Erlangen, Germany

Karl Förlinger

University of Tennessee
Knoxville, Tennessee

Al Geist

Oak Ridge National Laboratory
Oak Ridge, Tennessee

Michael Gerndt

Technische Universität München
Munich, Germany

Tom Goodale

Louisiana State University
Baton Rouge, Louisiana

Tobias Gradl

Friedrich–Alexander–Universität
Erlangen, Germany

William D. Gropp

Argonne National Laboratory
Argonne, Illinois

Robert Harkness

University of California, San Diego
San Diego, California

Albert Hartono

Ohio State University
Columbus, Ohio

Thomas C. Henderson

University of Utah
Salt Lake City, Utah

Bruce A. Hendrickson

Sandia National Laboratories
Albuquerque, New Mexico

Alfons G. Hoekstra

University of Amsterdam
Amsterdam, The Netherlands

Philip W. Jones

Los Alamos National Laboratory
Los Alamos, New Mexico

Laxmikant Kalé
University of Illinois
Urbana-Champaign, Illinois

Shoaib Kamil
Lawrence Berkeley
Berkeley, California

Cetin Kiris
NASA Ames Research Center
Moffett Field, California

Uwe Küster
University of Stuttgart
Stuttgart, Germany

Julien Langou
University of Colorado
Denver, Colorado

Hans Petter Langtangen
Simula Research Laboratory and
University of Oslo, Norway

Michael Lijewski
Lawrence Berkeley National
Laboratory
Berkeley, California

Anders Logg
Simula Research Laboratory and
University of Oslo, Norway

Justin Luitjens
University of Utah
Salt Lake City, Utah

Kamesh Madduri
Georgia Institute of Technology
Atlanta, Georgia

Kent-Andre Mardal
Simula Research Laboratory and
University of Oslo, Norway

Satoshi Matsuoka
Tokyo Institute of Technology
Tokyo, Japan

John M. May
Lawrence Livermore National
Laboratory
Livermore, California

Celso L. Mendes
University of Illinois
Urbana-Champaign, Illinois

Dieter an Mey
RWTH Aachen University
Germany

Tetsu Narumi
Keio University
Japan

Michael L. Norman
University of California, San Diego
San Diego, California

Boyana Norris
Argonne National Laboratory
Argonne, Illinois

Yousuke Ohno
Institute of Physical and Chemical
Research (RIKEN)
Kanagawa, Japan

Leonid Oliker
Lawrence Berkeley National
Laboratory
Berkeley, California

Brian O'Shea
Los Alamos of The National
Laboratory
Los Alamos, New Mexico

Christian D. Ott
University of Arizona
Tucson, Arizona

James C. Phillips
University of Illinois
Urbana-Champaign, Illinois

Simon Portegies Zwart
University of Amsterdam,
Amsterdam, The Netherlands

Thomas Radke
Albert-Einstein-Institut
Golm, Germany

Michael Resch
University of Stuttgart
Stuttgart, Germany

Daniel Reynolds
University of California, San Diego
San Diego, California

Ulrich R de
Friedrich–Alexander–Universit t
Erlangen, Germany

Samuel Sarholz
RWTH Aachen University
Germany

Erik Schnetter
Louisiana State University
Baton Rouge, Louisiana

Klaus Schulten
University of Illinois
Urbana-Champaign, Illinois

Edward Seidel
Louisiana State University
Baton Rouge, Louisiana

John Shalf
Lawrence Berkeley National
Laboratory
Berkeley, California

Bo-Wen Shen
NASA Goddard Space Flight Center
Greenbelt, Maryland

Ola Skavhaug
Simula Research Laboratory and
University of Oslo, Norway

Peter M.A. Sloot
University of Amsterdam
Amsterdam, The Netherlands

Erich Strohmaier
Lawrence Berkeley National
Laboratory
Berkeley, California

Makoto Taiji
Institute of Physical and Chemical
Research (RIKEN)
Kanagawa, Japan

Christian Terboven
RWTH Aachen University,
Germany

Mariana Vertenstein
National Center for Atmospheric
Research
Boulder, Colorado

Rick Wagner
University of California, San Diego
San Diego, California

Daniel Weber
University of Oklahoma
Norman, Oklahoma

James B. White, III
Oak Ridge National Laboratory
Oak Ridge, Tennessee

Terry Wilmarth
University of Illinois
Urbana-Champaign, Illinois

Symbols

Symbol Description

α	To solve the generator maintenance scheduling, in the past, several mathematical techniques have been applied.	$\theta\sqrt{abc}$	annealing and genetic algorithms have also been tested. This paper presents a survey of the literature
σ^2	These include integer programming, integer linear programming, dynamic programming, branch and bound etc.	ζ	over the past fifteen years in the generator
\sum	Several heuristic search algorithms have also been developed. In recent years expert systems,	∂	maintenance scheduling. The objective is to
abc	fuzzy approaches, simulated	sdf	present a clear picture of the available recent literature
		ewq	of the problem, the constraints and the other aspects of
		$bvcn$	the generator maintenance schedule.



Part I

**This is What a Part Would
Look Like**



1

Natural Language Processing

Author Name

CONTENTS

In this chapter we introduce modern NLP libraries, techniques and their applications. This chapter will focus on deep learning methods and less on computational linguistics that require nuanced knowledge of linguistics. We explore what it means to represent words and sequences of words with rich numeric representations that are better-suited toward modern computational tasks. We aim to capture some of these modern fine-tuned representations that are specially catered toward a semantic lexicon for medical language. We use these representations and aforementioned tools to showcase a modern reference implementation leveraging PyTorch, PyTorch Lightning and the Huggingface Transformers library. To wrap it all together, we walk through a complete example that highlights best practices that encourage reproducibility and allow for systematic iterative improvements.

This includes:

- An introduction into fundamental NLP concepts
- Bootstrapping techniques to iterate on a dataset in the low-resource setting
- Storing of a reference dataset in a publicly-accessible location
- Downloading, caching, loading, splitting, and preprocessing of the data
- Setting up of a cloud-based GPU workstation (?) (–this might be overkill for now, but keep if we can)
- VSCode (?)
- Monitoring the training run:
 - Logging and experiment tracking
 - Learning curves
 - Metrics

- Hyperparameter tuning, some tricks of the trade
- Offline evaluation and sanity checking

We will keep the discussion focused on tasks in epilepsy, using a running example of SUDEP prediction from electronic medical record (EMR) notes. Many of the concepts introduced here are very general and are straightforward translations to domains outside of SUDEP prediction, epilepsy, and even NLP.

1.0.1 Introduction to Natural Language Processing

Natural language processing (NLP) is a field of computer science that deals with the extraction, processing, and understanding of human language. It is known as the field of computer linguistics, and is a subfield of artificial intelligence. Common NLP tasks include sentence segmentation, tokenization, part-of-speech tagging, named-entity recognition, parsing, question answering, summarization and classification.

How can we teach a computer to perform these tasks? The first challenge is that computers, at their core, only understand numbers. So first we need to think about how words can be represented with numbers such that we can perform calculations on them. The numbers should allow the computer to assign meaning to words and their context with other words around it.

A simple way to achieve this is to define a vocabulary V of words, and assign each word a unique integer i . The word is then represented as a vector w of length $|V|$ with all zeros and a one at index i ¹.

```
have = [1, 0, 0, 0, 0, 0, ... 0]
a     = [0, 1, 0, 0, 0, 0, ... 0]
good = [0, 0, 1, 0, 0, 0, ... 0]
day  = [0, 0, 0, 1, 0, 0, ... 0]
...
```

We could now represent a sentence as the sum of the word vectors $S = \sum_j w_j$. This representation is suitable as input for any classification algorithm (such as logistic regression or a decision tree), to make a prediction about our target variable, e.g. whether the sentence is relevant to our epilepsy task. This simple representation fulfills our requirements but comes with some drawbacks:

- We implicitly assume that each word in the sentence is equally important.
- Each word is equally similar to every other word (e.g. by taking the euclidian distance between word vectors).
- The representation is invariant to reordering of the words.

¹This is also known as a one-hot encoding

The first point is a problem, because as we add more word vectors together, the sentence representation will converge to the global average and drown out any signal relevant to the specific sentence at hand. To address this we can instead write a weighted sum $S = \sum_j \lambda_j w_j$, where each word vector is weighted by its importance λ_j , and there are statistical methods we can use to compute an importance value². For example a word like *the* is very common and appears in many different contexts. It is unlikely that there is a lot of signal we can extract from it. On the other hand, a word like *epilepsy* will be much more rare and specific to a given task, and we would like to raise its importance. This has the net effect of allowing us to find a good representation for larger word sequences.

To address the second point, we will be touching on a very important concept, not just relevant in NLP, but also in the world of ML in general: Embeddings.

1.0.2 Embeddings

An embedding is a representation of a vector space that captures the similarity between the entities we are encoding, where entities can be words, images, e-commerce products, movies, houses, and many other data modalities. Scientific progress in the field of ML is often directly about finding algorithms that can learn high-quality robust embeddings in an efficient and scalable way, and the final performance of a given architecture is largely determined by the quality of the embeddings it learns.

To illustrate this concept, we will be talking briefly about one of the earliest algorithms that has been used to learn embeddings for words: *word2vec* [?], a technique pioneered by a team led by Thomas Mikolov. The idea behind *word2vec* is to set a word into context with the words surrounding it. In NLP, the study of techniques that determine the probability of a given sequence of words is known as language modeling.

For example the word *epilepsy* may appear in the context of words like *symptom*, *medication*, or *episode* and vice versa. However, we would not expect that *epilepsy* would appear in the context of words like *volleyball*.

We then translate this idea into a training task: Given the sum of the one-hot encoded word vectors of the words surrounding the target word, map it to the one-hot encoded vector of the target word.

This is a task that can be solved by a neural network with the following setup: Add a layer l_1 that maps from the input dimension (of size $|V|$) to an intermediate dimension of size h and then another layer l_2 that maps from the intermediate dimension to the output dimension (of size $|V|$). Here h is what is called embedding dimension and $h \ll |V|$. While h is generally a hyperparameter that can be optimized at a later stage, a reasonable choice is

²TF-IDF

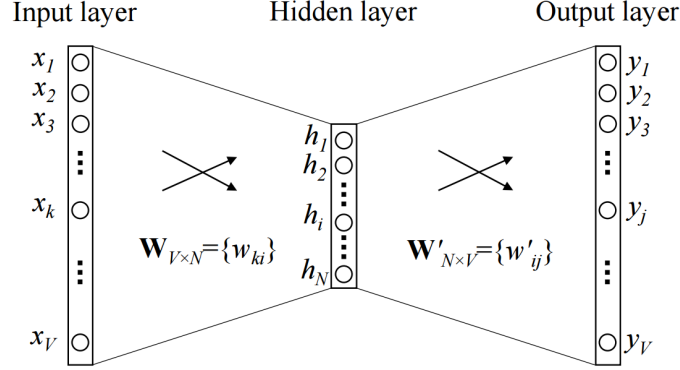


Figure 1: A simple CBOW model with only one word in the context

FIGURE 1.1
Word2vec

to use a value according to the rule of thumb:

$$h \approx 4\sqrt[4]{|V|} \quad (1.1)$$

For a vocabulary size $|V| = 30000$ this comes to a value of about 50. During training, the network will learn to encode the original sparse indicator vectors in a way that preserves the maximum amount of information necessary to predict the target word, while being forced to squeeze the information through the low dimensional bottleneck. An interesting side effect of this approach is that we can use part of the network to encode a given word in V by using for example the inverse of l_2 as a lookup table. The vectors that we obtain from this lookup table are called word embeddings and they have some very useful properties:

- Similar words tend to be close together (e.g. in terms of euclidian distance).
- Averaging word embeddings of a sentence will give us a more robust representation than our naive approach.
- We can perform some arithmetic, such as adding and subtracting the embedding vectors to traverse the space.

These properties also imply that the representations for two sentences will be close together if they are semantically similar, which will make the job of, say, a downstream classifier much more straightforward - it only has to slice the embedding space. The hard part of understanding the meaning of the sentence is already done. Since the training data can be generated from just raw text,

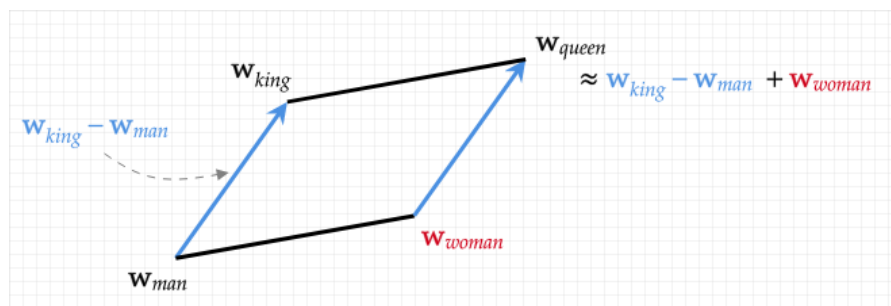


FIGURE 1.2
Word2vec

it is sufficient to train the embeddings once on a very large corpus³ and then reuse them, in effect giving us a headstart when building task-specific models. Instead of also having to learn what words mean, and some of the syntactic rules behind written language, we can get a head start on this and focus our computational resources on adapting this knowledge to a specific task.

This underlying assumption that similar neighboring words share a semantic similarity to center words is a powerful assumption, and what at the heart of many modern training objects for embeddings that are used in many downstream ML systems. We do not inherently inject a great deal of model bias, but rather rely on having a large enough training set that the word vectors will become meaningful representations with respect to one another. A corpus such as Common Crawl, which contains roughly 840 billion tokens of over two million unique words, provide such scale. There are other corpora used in these unsupervised processes, one such common candidate being a collection of all Wikipedia articles.

In practice, one training objective can be to use the distributed representation of a set of context words to guess a center word. This is known as the continuous bag of words (CBOW) model.

Consider a training document describing patients with epilepsy, where we are attempting to train a model with a predefined vocabulary V of the top For a vocabulary V , where we have chosen the top $|V|$ most common words in English to learn. The document being chosen in this training iteration may contain the following excerpt:

```
... recommended administering Lamictal to treat their
focal seizures, the patient continued experiencing
symptoms though reported a 50% reduction in ...
```

In training our word embeddings, we will randomly mask one of the words,

³GloVe

in this case *patient*. Our model makes use of an additional parameter, a context window size c , that helps to provide necessary context. If we choose $c = 8$, our training sample would become

```
... lamictal to treat their focal seizures , the [TARGET]
continued experiencing symptoms though reported a 50 % ...
```

For a given center word at index i , we consider previous words at indexes $i - 1, i - 2 \dots i - c$, and all subsequent words at indexes $i + 1, i + 2 \dots i + c$. In our word sequence S , we may then choose to sum the word vectors of each of these context words. Our context input then becomes:

$$\sum_{\substack{j=-c \\ j \neq 0}}^c w_{i+j} \quad (1.2)$$

Our task then becomes to predict a probability distribution over the size of our vocabulary $|V|$, where the target is a vector of all 0's with a 1 at the index in the vocabulary for the word *patient*

There are a few additional takeaways from this example. Note that we will address many of these in the sections to follow.

- The training process will consider this text probably many times, and mask out different words in future iterations.
- Our processed sentence is slightly modified from our original input. Notice that the number and percent symbol are split, and that Lamictal is lowercased.
- The word *patient* can also be an adjective. Our *patient* embedding will also capture this meaning in our model, dependent upon how often that meaning is used.
- If our context window is too small, we lose information. If it is too large, our signal becomes fuzzy as we converge to the global mean vector.

What if we choose to use only a center word, and have our model predict all of the surrounding context words? This is a common training technique for embeddings as well, and is known as the Skip-Gram model.

1.0.2.1 Character and Subword Embeddings

This is placeholder text

1.0.2.2 Vocabularies and Medical Embeddings

This is placeholder text

1.0.3 Transformers

Until recently building strong NLP systems required deep understanding of language and its structure, as well as large amounts of data, even when relying on pre-trained word embeddings. The resulting systems were often still brittle.

In recent years some surprisingly powerful architectures and pre-training regimes have emerged that build on the concepts we have discussed here, most notably under the name of *Transformers*⁴, that address the last bullet point of our representation’s shortcomings.

The details of the inner workings are beyond the scope of this tutorial, for now it suffices to know that they are able to learn robust sentence embeddings with a fine grained understanding of syntactic structure, semantics, and general knowledge of the world. In fact, they can be considered near the level of a human that just knows the English language (or other languages), along with a broad array of factual knowledge about the world (think Wikipedia). Just like a human, they can be taught specialty knowledge that is relevant to a given task, in ML we would say we *finetune* them. This paradigm shift has led to a step function improvement in terms of performance and sample efficiency in a wide variety of NLP tasks where Transformers are the defacto state of the art.

In the following sections, we will discuss a broad set of problem statements Transformers are suitable for, along with a set of techniques a practitioner can employ to arrive at robust solutions, even with limited data.

1.0.4 NLP Tasks in Epilepsy

This is a placeholder title. A couple pages here might do nicely as example studies for some important techniques for epileptologists. Now that we have an understanding of some of the basics of computational representations of language data, let us take a look at some of the modern use cases for NLP in industry, catered toward epileptologists and other medical practitioners.

1.0.4.1 Text Preprocessing and Tokenization

“nobreak

1.0.4.2 Named Entity Recognition and PHI

This is a working title. NER is an important problem in medicine, for PII if nothing else, and I figure we could go through some basics of using it with Spacy, NLTK, or Stanford NLP here. Stopwords would be worth mentioning.

⁴BERT, GPT

1.0.5 Mining Unstructured Text and Clinical Documentation

This is placeholder text. Talk about clinical documentation and EHR notes, and other forms of unstructured text, being used for downstream classification tasks. Then, give some examples of them, and what these people did, and possibly how transformers could probably improve this work greatly.

1.0.5.1 Clinical Trial Matching and Eligibility Prescreening

This is a placeholder text. <https://www.nature.com/articles/d41586-019-02871-3>

1.0.5.2 Surgical Candidacy

This is placeholder text

1.0.5.3 Avoiding Epilepsy Misdiagnosis

This is placeholder text. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4419916/>

1.0.6 Full Walkthrough - SUDEP Detection

This is placeholder text