# Title line 1
# Title line 2

**A thesis submitted in partial fulfilment**

**of the requirement for the degree of Doctor of Philosophy**

# Name M. Lastname

# July 2011

# Cardiff University
# School of Computer Science & Informatics

## Declaration

This work has not previously been accepted in substance for any degree and is not concurrently submitted in candidature for any degree.

Signed . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . (candidate)

Date    . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

## Statement 1

This thesis is being submitted in partial fulfillment of the requirements for the degree of PhD.

Signed . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . (candidate)

Date    . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

## Statement 2

This thesis is the result of my own independent work/investigation, except where otherwise stated. Other sources are acknowledged by explicit references.

Signed . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . (candidate)

Date    . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

## Statement 3

I hereby give consent for my thesis, if accepted, to be available for photocopying and for inter-library loan, and for the title and summary to be made available to outside organisations.

Signed . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . (candidate)

Date    . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .

**To People you care
for their patience and support.**

# Abstract

We produce interpretable representations, and demonstrate their applicability in interpretable classifiers. Our approach is model-agnostic, given a similarity-based representation, we are able to produce a representation in terms of domain knowledge. We evaluate the interpretability of our representation and provide examples of interpretable classifiers with our representation.

# Acknowledgements

# Contents

# List of Publications

The work introduced in this thesis is based on the following publications.

- 

-

# List of Figures

# List of Tables

# List of Algorithms

# List of Acronyms

**ML** Machine Learning

**NLP** Natural Language Processing

**NDCG** Normalized Discounted Cumulative Gain

## 0.0.1 Definitions

**Domain** Where the data was originally sourced from $DOM^I MDB$, e.g. IMDB movie reviews.

**Word** A string of alphanumeric characters that originated from text in the domain $DOM_w$, e.g. the $w = "Horror"$ from a domain of IMDB movie reviews $DOM^I MDB$.

$w$

**Corpus of Documents** A unique group of words, e.g. a review from a domain of IMDB movie reviews $DOM_I MDB$.

$C_d w$

**Document** A document of words

$d_w$

**Vector Space** A representation composed of vectors.

$S_v$

**Semantic Space** A representation where spatial relationships between vectors correspond to semantic relationships.

$S_v$

**Word frequency** The frequency of a word $wf$ for its document $D_wf$.

$wf$

**Bag-Of-Words** a matrix BOW of documents $BOW_D$ where each document is composed of unordered frequencies of words $D = [wf_1, ..., wf_n]$. and Conceptual Space we obtain a representation of entities composed of properties. Then, we cover the additional methods we propose to improve this process.

$BOW_d$

**Bag-Of-Words PPMI**

**Feature** A feature is a distinct useful aspect of the domain, corresponding to a numerical value.

$R_f$

**Hyper-plane** The hyper-plane for a word

$H_w$

**Direction vector** The orthogonal direction to a hyper plane that separates a word in a vector space.

$D_w$

**Cluster label** A cluster of words that describe a property.

$C_w$

**Cluster direction** The averaged directions of all words in the label.

$D_C$

**Feature rankings** The rankings induced from a feature direction.

$R_D C$

*Chapter 1*

# Introduction

## 1.1 Motivation

With the rise of services on the web that enable large-scale user-generation of text data, e.g. Social Media sites (Facebook, Twitter), Review sites (IMDB, Rotten Tomatoes, Amazon) and content-aggregation sites (Reddit, Tumblr), the internet has become largely populated by text posts that are related to some specific, niche topic within a domain. For example, a review on Amazon for a product is specially tailored text for that product within the domain of Amazon reviews. Taken from a closer lens, we could even argue that each review-type has its own domain, e.g. Product reviews, Food reviews, Movie reviews. However, the text posts themselves are largely unstructured semantically. Humans can have an intuitive understanding of the semantics that are present in unstructured text, but machines do not.

One task of Natural Language Processing is to obtain this semantic understanding from text by obtaining a machine-readable representation that contains domain knowledge. A basic approach to obtain a representation of this text is to represent entities (e.g. reviews, text-posts) by the frequency of their words, see 1.1.

Below, we show a review with its associated properties labelled.

We can understand these properties to have a degree to which they apply, for example the size of the clothing might be "XXL", "XL", "L", "M" or "S", or the quality may be "Very good", "Good", "Ok", "Bad" or "Very bad". For the former, we may rely on the metadata available from the site itself, but for the latter the way to obtain this information is less clear. Although we may infer that the rating has some indication of these properties, it does not describe the properties or the degree to which the review refers to them. This kind of information is valuable

**Figure 1.1: Bag-of-words**



**Figure 1.2: Example properties**

for making sense of the world of unstructured text, and has broad applications, e.g. The most immediate example is perhaps that they allow for a natural way to implement critique-based recommendation systems, where users can specify how their desired result should relate to a given set of suggestions [**?**]. For instance, [**?**] propose a movie recommendation system in which the user can specify that they want to see suggestions for movies that are "similar to

this one, but scarier". If the property of being scary is adequately modelled as a direction in a semantic space of movies, such critiques can be addressed in a straightforward way. Similarly, in [**?**] a system was developed that can find "shoes like these but shinier", based on a semantic space representation that was derived from visual features. Semantic search systems can use such directions to interpret queries involving gradual and possibly ill-defined features, such as "*popular* holiday destinations in Europe" [**?**]. While features such as popularity are typically not encoded in traditional knowledge bases, they can often be represented as semantic space directions.

### 1.1.1 Directions

However, manually labelling these properties and the degrees to which entities (e.g. reviews, text-posts) have them is extremely time-consuming.

A potentially ideal system would be as follows: We collect large amounts of unstructured text data, separated into domains, and obtain the properties of each domain from this data, and rank entities on the degree to which they have these properties. In this way, properties would be understood on a scale built from the domain directly, so that each domain has its own meanings for words according to their own idiosyncrasies. As the process does not require any manual labelling the quality of these properties could be improved simply by obtaining more data. Further, as we are learning from unstructured data, not only would this allow us to understand the data in terms of what we know, but it would also introduce us to new ideas that we may not have previously understood. This kind of representation also has value in application to Machine Learning tasks. If we can separate the semantics of the space linearly into properties, we are able to learn simple linear classifiers that perform well.

Simple linear classifiers built from a representation composed of rankings on properties have an additional benefit of being more understandable.

## 1.2   Interpretability

Most successful approaches in recent times, like vector-spaces, word-vectors, and others, rely on the distributional model of semantics. This model relies on encoding unstructured text e.g. of a movie review, as a vector, where each dimension corresponds to how frequent each word is, we are able to calculate how similar the entities are, e.g. we know that if two movies have a similar distribution of words in their reviews, like frequent use of the word 'scary', or 'horror', then they would have a higher similarity value. These models, also known as 'semantic spaces' encode this similarity information spatially.

Semantic relationships can be obtained from semantic spaces.

applications/need for good interpretability:

- Safety

- Troubleshooting, bug fixing, model improvement

- Knowledge learning

- EU's "Right to explanation"

- Discrimination

properties of an interpretable classifier:

- Complexity: 'the magic number is seven plus or minus two' [5] also has many positive effects for its users, like lower response times [4, 2], better question answering and confidence for logical problem questions [2] and higher satisfaction [4].

- Transparancy:

- Explainability:

- Generalizability:

Properties, entities, the benefits and application of a representation formed of these

Basic introduction to directions, explanation of the utility and application of our approach

## 1.3    Thesis Overview / Contributions

In 3, we focus on further experimenting with one relationship that was formalized in [1]: a ranking of entities on properties. In particular, we use this method of building a representation of entities as a way to convert a vector space into an interpretable representation, for use in an interpretable classifier. The reason that we chose this representation to expand on is because by representing each entity $e$ with a vector $v$ that corresponds to a ranking $r$, the meaning of each dimension is distinct, and we are able to find labels composed of clusters of words for these dimensions. Here, we make the distinction between a property and a word, a property is a natural property of the space that exists in terms of a ranking of entities, and words are the labels we use to describe this property.

*Chapter 2*

# Background

## 2.1 Text Representations

Need to write about the concept of salient features of a domain here.

### 2.1.1 Bag-of-words

We begin by processing an unstructured text corpus, composed of documents $C_D$. We then remove all punctuation, convert any accented characters to non-accented characters, and lower-case the documents to obtain word tokens for each document $D_W$. From here, we can assume that any $W \approx W$ will now $W = W$, if a word varied in format but not alphanumeric characters.

Then, we count the occurrences of each word

- Frequency

- Tf-idf

- PPMI

## 2.2 Text classification

### 2.2.1 Decision Trees

- Explanation of what decision trees are

- Explanation that they may not perform well on sparse information

- Max features

- Criterion

- CART decision trees versus others

### 2.2.2 Support Vector Machines

- Performance increase for support vector machines on sparse data, balancing, etc

- C parameters, gamma parameters

### 2.2.3 Neural Networks

- Difference between SVM and Nnet

### 2.2.4 Semantic Spaces

Bag-Of-Words representations of text result in large sparse vectors for each document,

**How do vector spaces represent semantics? Why do we use them to represent semantics?**

Distributional representations of semantics, known as 'semantic spaces' are well-recognized for their ability to represent semantic information spatially. These representations have been widely adopted for Natural Language Processing (NLP) tasks thanks to their ability to represent complex information in a dense representation. In particular, entity-embeddings have been applied to represent items in recommender systems [**?**, **?**, **?**], to represent entities in semantic search engines [**?**, **?**], or to represent examples in classification tasks [**?**].

Vector spaces are a popular way to represent unstructured text data, and have been broadly applied to and transformed by supervised approaches. They vary in method, producing structure from Cosine Similarity, Matrix Factorization, Word-Vectors/Doc2Vec, etc. They also vary in how they linearly separate entities. However, their commonality is that they are able to represent

semantic relationships spatially. See Section 2.2.4 This brings up an essential point: When using a semantic space, are we taking advantage of relationships that are discriminative or incorrect? The danger of relying on these spaces and the models that use them has greatly affected their adoption in critical application areas like medicine, and has raised legal concerns about their application in e.g. determining if someone is suitable for a loan.

See Section 2.2.4

- Word-vectors

### 2.2.5   Document Representations

**LSA**

Principal Component Analysis is a dimensionality reduction method that results in dimensions ordered by importance. Starting with a large data matrix, e.g. our TF-IDF values from before, we first find the covariance matrix for these values. Then, from this covariance matrix we obtain the eigenvalues. We can then linearly transform the old data in-terms of this covariance matrix to obtain a new space of size equal to an arbitrary value smaller than our matrix.

- PCA

- MDS

## 2.3   Interpretable Representations

a. NNSE b. compositional c. 2007 paper as wikipedia similarities d. Topic models e. Infogan, etc

[**?**] Sparse PCA (Why not compare lol)

Vector space models typically use a form of matrix factorization to obtain low-dimensional document representations. By far the most common approach is to use Singular Value Decomposition [**?**], although other approaches have been advocated as well. Instead of matrix factorization,

another possible strategy is to use a neural network or least squares optimization approach. This is commonly used for generating word embeddings [**?**, **?**], but can similarly be used to learn representations of (entities that are described using) text documents [**?**, **?**, **?**]. Compared to topic models, such approaches have the advantage that various forms of domain-specific structured knowledge can easily be taken into account. Some authors have also proposed hybrid models, which combine topic models and vector space models. For example, the Gaussian LDA model represents topics as multivariate Gaussian distributions over a word embedding [**?**]. Beyond document representation, topic models have also been used to improve word embedding models, by learning a different vector for each topic-word combination [**?**].

The most commonly used representations for text classification are bag-of-words representations, topic models, and vector space models. Bag-of-words representations are interpretable in principle, but because the considered vocabularies typically contain tens (or hundreds) of thousands of words, the resulting learned models are nonetheless difficult to inspect and understand. Topic models and vector space models are two alternative approaches for generating low-dimensional document representations.

## 2.3.1 Word Vectors

*Chapter 3*

# Converting Vector Spaces into Interpretable Representations

## 3.1 Introduction

The ever more pervasive digital infrastructure that supports our lives has resulted in many opportunities to obtain data and models to make sense of that data. Semantic Spaces that encode semantic relationships between documents spatially have recently achieved strong results on tasks like X, Y, Z. These neural-network learned representations make use of a variety of new information like grammatical structure, word-context and even image data. Further, as domains become more entrenched in the digital world, the need for models in safety critical domains like medicine or legal domains like credit evaluation have increased the need for producing interpretable models, as well as interpretable representations. However, the dimensions of a semantic space do not correspond to human understandable features, and standard approaches to interpretable text representations do not match the performance of these methods. Ideally, we would obtain a representation that makes use of the rich semantic relationships from a high-performing semantic space, but also has dimensions corresponding to interpretable features. To this end, we aim to introduce in this chapter a methodology to linearly transform a semantic space using just its associated bag-of-words as input into an interpretable representation, and demonstrate the applicability of this interpretable representation to simple interpretable classifiers.

There are many types of semantic relationships in a semantic space. For our work, the representation is composed of rankings of documents on semantic directions in the space, in particular where those directions correspond to features. We show an example of the kind-of directions we use to obtain our representation in 3.1. Directions from domain-specific semantic spaces have

**Figure 3.1: An example in a toy domain of shapes.**

been used previously in a variety of ways, For instance, [?] found that features of countries, such as their GDP, fertility rate or even level of $CO_2$ emissions, can be predicted from word embeddings using a linear regression model. Similarly, in [?] directions in word embeddings were found that correspond to adjectival scales (e.g. bad $<$ okay $<$ good $<$ excellent) while [?] found directions indicating lexical features such as the frequency of occurrence and polarity of words.

Derrac [1] introduced an unsupervised method to go from a semantic space and its associated bag-of-words to a representation where each dimension is a ranking of documents on a feature of the domain. For example, in the domain of movie reviews genres would be a feature, and the dimension would have a numeric value for each document corresponding to the degree it is a particular genre. The contribution of this Chapter is an analysis and experimentation on the

and rep.png

| Movie name | Properties and associated ranking |
|---|---|
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |
| Titanic | Romantic: 95%, Dramatic: 80%, Boat: 90%, Scary: 40% |

**Figure 3.2: Example movies and selected associated dimensions, chosen according to their relevance to the genre task..**

quality of these features applied to document classification. The main insight from our work is that these interpretable features do not suffer a performance drop in a non-linear classifier compared to the original representation, and can outperform the original representation and a baseline interpretable representation in a linear classifier. In addition, we find that if a dimension ranks documents on a feature relevant to the task, it can be competitive with more complex models using a single decision tree node. We show an example of the representation from a domain of IMDB movie reviews in 3.2.

This chapter continues as follows: We begin by describing related work, then explain the method, making explicit the variations we have introduced for our experimental work. We follow this with the results of our experiments accompanied by qualitative examples and explanations, and finish with a conclusion on the benefits and limitations of this approach.

These relationships have been expanded on, for example [6] found that "equivalent relations tended to correspond to parallel vector differences" [3], and [3], found that by decomposing representations into orthogonal semantic and syntactic subspaces they were able to produce substantial improvements on various tasks. Additionally, they have also been found to hold inherent gender bias [?] as word distances between gendered words (e.g. male, female, she, her) and occupational words e.g. (nurse, programmer) were correlated to the percentage of occupation that gender had for that role in different time periods.

## 3.2 Method

This section details the methodology to go from a Bag-Of-Words (BOW) 2.1.1 and Semantic Space 2.2.4, to rankings of documents on features of the domain, e.g. In a domain of IMDB movie reviews, where a document is composed of all of its reviews, a movie would be ranked on features like $Scary, Horror, Bloody$ and $Romantic, Love, Cute$, ideally with as many rankings as salient features of the domain.

### 3.2.1 Obtaining Directions and Rankings From Words

In this section we show how to obtain directions for words, and explain how to obtain document representations by ranking documents on these directions. For this step, we do not expect all words to be features of the domain. In the next sections, we aim to filter these words to obtain salient features.

**Obtaining directions for each word** For each word $w$, a Support Vector Machine (See Section 2.2.2) classifier is trained on the binary Bag-Of-Words representation of that word, where words are labelled as positive if they occurred more than once $w_f >= 1$ and negative otherwise. Although the separation of documents is binary, we can expect that the degree to which they are classified as the word varies. For example in a space constructed from frequency vectors, we can expect that the documents which contain the word more frequently would be further away from the hyper-plane in the positive direction. Following this, we can consider the vector $v_w$ perpendicular to the hyperplane as the direction that models documents from least relevant at the distance furthest from the hyperplane on the negative side to most relevant for the word $w$ at the distance furthest from the hyperplane at the positive side. We show an example of this in the toy domain in Figure 3.3.

**Ranking documents on directions** Once we have obtained a direction vector for each word $v_w$ the next step is to quantify the degree to which each document has that word, by obtaining a value that corresponds to how far-up it is on the direction vector. These are our rankings of documents on words, if $p_d$ is the representation of a document in the given vector space as a point then we can think of the dot product between the hyper-plane and the document vector $H_w \cdot p_d$ as the ranking $r_d w$ of the document $d$ for the word $w$, and in particular, we take $r_d 1 < r_d 2$

**Figure 3.3: An example of a hyper-plane and its orthogonal direction in a toy domain of shapes. Green shapes are positive examples and red shapes are negative examples, but despite the problem being binary those closest to the hyper-plane are less defined than those further away, resulting in the orthogonal vector being a direction..**

to mean that $d_2$ has the property labelled with the word $w$ to a greater extent than $e_1$. Below, we show some examples of features and documents ranked on them for different domains.

### 3.2.2 Filtering Words

With the rankings $R_r$, we could create a representation of each document $d$, composed of $w_n$ dimensions, where each dimension is a ranking of the document $d$ on that word $r_d w$. However,

many of the words are not spatially important enough in the representation to result in a quality ranking - they are not salient features. In this section, we aim to filter the words that are not separable, we evaluate them using a scoring metric, and remove the words that are not sufficiently well scored. We use three different metrics:

**Classification accuracy**. Evaluating the quality in terms of the accuracy of the SVM classifier: if this classifier is sufficiently accurate, it must mean that whether word $w$ relates to document $d$ (i.e. whether it is used in the description of $d$) is important enough to affect the semantic space representation of $d$. In such a case, it seems reasonable to assume that $w$ describes a salient property for the given domain.

**Cohen's Kappa**. One-kind of feature we find in these domains are binary labels of documents, for example a movie either is or isn't a movie with "Gore". We can expect that the more salient a binary feature, the more linearly separable it will be in the space. One problem with accuracy as a scoring function is that these classification problems are often very imbalanced. In particular, for very rare words, a high accuracy might not necessarily imply that the corresponding direction is accurate. For this reason, [1] proposed to use Cohen's Kappa score instead. In our experiments, however, we found that accuracy sometimes yields better results, so we retain Kappa as an alternative metric.

**Normalized Discounted Cumulative Gain** This is a standard metric in information retrieval which evaluates the quality of a ranking w.r.t. some given relevance scores [**?**]. In our case, the rankings $r_d$ of the document $d$ are those induced by the dot products $v_w \cdot d$ and the relevance scores are determined by the Pointwise Positive Mutual Information (PPMI) score $ppmi(w, d)$, of the word $w$ in the BoW representation of entity $d$ where $ppmi(w, d) = \max\left(0, \log\left(\frac{p_{wd}}{p_{w*} \cdot p_{*d}}\right)\right)$, and

$$p_{wd} = \frac{n(w, d)}{\sum_{w'} \sum_{d'} n(w', d')}$$

where $n(w, d)$ is the number of occurrences of $w$ in the BoW representation of object $d$, $p_{w*} = \sum_{e'} p_{wd'}$ and $p_{*d} = \sum_{w'} p_{w'd}$.

By scoring the words on these features, we can apply a simple cut-off (e.g. the top 2000 scored words) to obtain the most salient words. Ideally, this cut-off would be at the point where the words stop corresponding to salient features. However, it is difficult to determine this. In

principle, we may expect that accuracy and Kappa are best suited for binary features, as they rely on a hard separation in the space between objects that have the word in their BoW representation and those that do not, while NDCG should be better suited for gradual features. In practice, however, we could not find such a clear pattern in the differences between the words chosen by these metrics despite often finding different words. In Table **??**, we show examples of the differences between the largest differences between the scoring methods.

**Clustering Direction Vectors**

If we consider two directions, "Blood" and "Gore", we can understand both of these to be approximating a similar feature of movies, as they both relate to how much blood a movie contains. Because of this, we can expect their directions to be very similar to each other. This is the first idea behind clustering these directions, if we average these directions together we can obtain a direction inbetween them that is a balance between documents that used the word 'Bloody' to describe the blood and the word 'Gore'. To expand on this, some entities would have the property of being bloody films, but did not necessarily use the term gore in their reviews, same as some entities having the property but using the term gore not bloody, we can understand that this new hyper plane and associated direction more accurately represents the property of a bloody film more than either of the terms individually. By extending this to a clustering method, we can find similar abstract features by ensuring that all similar directions are clustered together.

The word direction for "beautiful" can be nebulous to the interpreter, as it is not clear what it means for a movie to be ranked highly on 'beautiful'. Considering this, clustering provides another advantage, once we cluster the terms to find the property ("beautiful", "cinematography" "shots") we are given context for the word and more easily intuit the feature, in this case it is a feature about how well the movie was directed.

We approach clustering the directions with a variety of methods:

**K-Means** K-Means is a clustering algorithm that starts with determining the amount of clusters, $K$. To begin, $K$ centroids $c$ are randomly placed into the space. Then, the distance between each point $p$ and centroid $c$ (in our case, points are determined by rankings) is calculated. Each point $p$ is then assigned to its closest centroid $c$. Then, the centroids are recomputed to be the

mean of their assigned points. This process starting with the distance calculation is repeated until the points assigned to the centroids do not change.

**Derrac's K-Means Variation** This is the clustering method used in the previous work [**?**]. As input to the clustering algorithm, we consider the $N$ best-scoring candidate feature directions $v_w$, where $N$ is a hyperparameter. The main idea underlying their approach is to select the cluster centers such that (i) they are among the top-scoring candidate feature directions, and (ii) are as close to being orthogonal to each other as possible.

The output of this step is a set of clusters $C_1, ..., C_K$, where we will identify each cluster $C_j$ with a set of words. We will furthermore write $v_{C_j}$ to denote the centroid of the directions corresponding to the words in the cluster $C_j$, which can be computed as $v_{C_j} = \frac{1}{|C_j|} \sum_{w_l \in C_j} v_l$ provided that the vectors $v_w$ are all normalized. These centroids $v_{C_1}, ..., v_{C_k}$ are the feature directions that are identified by our method.

We choose our first cluster centroid by taking the top-scoring direction for its associated metric. Then, we select centroids until we have reached the desired amount by taking the maximum of the summed absolute cosine similarity of all current centroids, in other words taking the most dissimilar direction to all of the current directions. Once we have selected the centroids, for each remaining direction we find the centroid it is most similar to, and the centroid is updated once the direction has been added.

## 3.3    Quantitative Results

### 3.3.1    Datasets

We use five different domains:

Newsgroups, originally containing 18,846 documents, is preprocessed using sklearn to remove headers, footers and quotes. Then, empty and duplicate documents are removed, resulting in 18302 documents. The vocabulary size is 141,321. The data is not shuffled. After filtering out terms that did not occur in at least two documents, we ended up with a vocabulary of size 51,064.

Sentiment is dataset where documents are reviews, containing 50,000 documents with a vocabulary size of 78588. After removing terms that did not occur in at least two documents, we ended up with a vocab of size 55384. Notably, this means that we removed all terms that did not occur in two documents for the sentiment, and in two documents for the newsgroups, and newsgroups began with a larger vocabulary than sentiment, but the ending vocabularies were about the same. This means that the terms in the newsgroups were more sparse than sentiment. In other words, newsgroups contained many terms that were not relevant to a majority of the documents. This is unsurprising, as it is a collection of 20 different newsgroups, rather than one single domain.

Reuters is a dataset of reuters news wires, originally containing 10788 documents. After removing empty and duplicate documents, we end-up with 10655 documents. It originally contained 90 classes, but as they were extremely unbalanced we removed all classes that did not have at least 100 positive instances, resulting in 21 classes. All other classes in other domains meet this threshold. The original vocabulary size is 51,0001, and after removing all words that do not occur in at least two documents, the vocabulary size is 22542.

Placetypes is a data-set of flickr tags, taken from the previous work [1]. It originally has a vocabulary size of 746,527 and 1383 documents. This is a very large vocabulary size to document ratio. The end vocabulary for this space was 100,000, which we used as a hard limit. This is roughly equivalent to removing all documents that would not be in at least 6 documents.

Movies is a dataset where each document is a movie represented by all of its reviews concatenated across a number of sources. It starts off with a vocabulary size of 551,080 and a document size of 15,000. However, after investigating the data made available by the authors, we found that there were a number of duplicate documents. After removing these duplicate documents, we end-up with 13978 documents. In the same way as the movies, we limit the vocabulary size at 100,000.

For all of these datasets, we split them into a 2/3 training data, 1/3 test data split. We additionally remove the end 20% of the training data and use that as development data for our hyper-parameters, which is then not used for the final models verified using test data.

### 3.3.2   Evaluation Method

We primarily examine the effectiveness of a representation on its ability to perform in low-depth Decision Trees, specifically CART Decision Trees (See Background Section 2.2.1) with a limited depth of one, two and three. We enter this evaluation with a few assumptions: A good interpretable representation disentangles salient domain knowledge into its dimensions, and natural domain tasks (e.g. classifying genres of movies using their reviews) can be evaluated effectively using that salient domain knowledge. Put another way, if the space is representing domain knowledge well we can expect that the space is linearly separable for key semantics of the domain. In spatial terms, a representation will be capable of being linearly transformed by our method into these distinct relevant concepts if semantically distinct entities are spatially separated, and semantically similar entities are close together.

If we only wanted to evaluate the quality of the representation, we could use Linear SVM's to find the hyper-planes that effectively separate these spatial representations for the class. However, the representations that encode this spatial information are not interpretable, so a linear classifier although able to separate the documents that contain the class and do not contain them will not be interpretable either. It is our main interest to evaluate how well a representation encodes these key semantics while also being restricted by the requirement to be disentangled into words or clusters, in other words how well it represents the information while also being interpretable.

Given these assumptions, low-depth decision trees can give an estimation of how good an interpretable representation is. If the representation cannot perform for a class at a one-depth tree, then it is not disentangled such that it contains a single salient dimension that effectively evaluates a class. If a representation cannot perform well on two-depth trees, then the representation is not disentangled into three concepts that can sufficiently determine that class, and if a representation cannot perform well on three-depth trees, it has not disentangled the representation such that there are nine relevant concepts that are relevant to that class. To see what these different trees look like see figure 3.4. A comparison to put this in better perspective is to an unbounded tree. Unbounded trees select a large amount of dimensions in order to achieve a performance difference on development data, but when applied to test data the models do not generalize well. This is because they overfit, rather than using the key semantics of the space to

**Figure 3.4: This figure shows an example tree from one of our classifiers. Here, we can see that the model increases in complexity as it increases in depth. In this case, we end-up getting better F-score with just a depth-one tree, as the tree begins to overfit at depth three. .**

classify.

We look primarily at the F1-score to determine if a classifier is good or not. This is because many of the classes are unbalanced (See above in Section 3.3.1 for exactly how unbalanced) so accuracy is not a good metric, as high accuracy could be achieved by predicting only zeros. All of the results shown in this section are the end-product of a two-part hyper-parameter optimization. Each Decision Tree has its own set of hyper-parameters that are optimized as does each representation-type. These are the models trained on the training data and scored on the test data, with the highest performing in terms of F1-score parameters from hyper-parameter optimization on the development data. For ease of comparison, we provide some results with SVM's and unbounded Decision Trees, as well as a baseline Topic Model, which we use as a reference for a standard interpretable representation. Below, we list the parameters that we optimize for each of these model types:

**Linear Support Vector Machines (SVM's)**: **Topic Models:CART Decision Trees** :

**Multi-Dimensional Scaling (MDS)**: **Principal Component Analysis (PCA)**: **Doc2Vec (D2V)**: **Average Word Vectors (AWV)**:

When obtaining the single word directions, we take all of the baseline representations and vocabularies, and filter the vocabularies according to a hyper-parameter that we tune. As the

**Figure 3.5: A conceptual space of movies, where regions correspond to properties and entities are points..**

doc2vec has already been hyper-parameter optimized, we use the optimal doc2vec space that scored the highest for its class on a Linear SVM, rather than tuning the entire process around the doc2vecs vectors. So for example, when we are evaluating the Keywords task for the movies, we would obtain directions from the doc2vec space that performed best for a linear SVM on the Keywords task following the previous experiments.

The parameters were We are not able to obtain an MDS space for sentiment or doc2vec spaces for placetypes/movies.

For example, a good vector space in the domain of movies constructed from IMDB movie reviews should contain a natural separation of entities into genres, where Horror movies are spatially distant from Romance movies, and movies that are Romantic Horrors would be somewhere inbetween. We can see an example in Figure 3.5. For a Bag-Of-Words, we can expect similar entities to have similarly scoring terms **??**.

We obtain results for the rankings induced from these word directions on Decision Tree's limited to a depth of 3 in-order to select the best parameters when using directions for each class. The parameters that we want to determine are the type of Semantic Space, the size of the space, the frequency threshold and the score threshold. To do so, for each space-type of each size, we use a grid search to find the best frequency and score cut-offs for that sized space-type. Then, we select from these space-types and sizes the best performing one. We can understand there to be a balance between finding words which are useful for creating salient features in our clustering step without including too many words which do not. As our clustering methods are unsupervised, it is important that we try and limit the amount of junk being entered into them,

| | Top PPMI scoring terms |
|---|---|
| Example Horror Entity | Term term term term term term term term term term term term term term |
| Similar Horror Entity | Term term term term term term term term term term term term term term |
| Somewhere Inbetween Entity | Term term term term term term term term term term term term term term |
| Romance Movie | Term term term term term term term term term term term term term term |
| Similar Romance movie | Term term term term term term term term term term term term term term |

**Table 3.1: Two of the following entities: Those classified as horror, those classified as horror and romance, and those classified as romance with their associated highest value PPMI terms. We show the highest positive instances here as the representation is sparse, even though we can also expect the terms that are low scoring to be similar too..**

despite the classifiers that use these directions typically being able to filter out those directions which are not suitable to the class. Additionally, as the vocabulary size varies from dataset to dataset, the threshold will naturally be different for each one.

These results allow us to choose for each class, the best Semantic Space and Scoring-type for that class. Next, we test single directions, attempting to find a good amount of directions to cluster and not including words which may hamper the unsupervised classification, as well as the best space-type for each domain. We found that generally, X was the best space and as expected classifiers performed better with more data, so we use 20000 as our frequency cutoff and 2000 as our score cutoff.

We continue with the optimal space and score-type chosen by our single direction experiments, and use the same frequency and score thresholds as before. We then experiment with two different clustering algorithms: Derrac and K-Means. As these algorithms select centroids from the top-scoring directions or randomly, we can expect that some clusters may not be salient features of the space. This is because top-scoring directions, e.g. for accuracy could simply infrequent terms that do not have much meaning, and these infrequehnt terms could also be randomly selected. We could use grid-search on the frequency and score cutoffs when obtaining these results in order to avoid terms that may disrupt existing clusters or form cluster centers that are not salient features of the space, but we chose a more standardized process that would rely on the parameters of the clustering algorithms and the ability of the classifiers to filter out clusters that are not informative, so as to not make a time-costly grid search a necessary part of

the process.

With that in mind, we use three clustering algorithms.

Mini batch K-means, implemented by scikit-learn [1], introduced by [**?**] and kmeans++ to initialize [**?**]

### 3.3.3 Summary of all Results

To begin, we compare the original dimensions of the space, the rankings on single words, the rankings on cluster directions, a bag-of-words of PPMI scores and

---

[1]https://scikit-learn.org/stable/modules/generated/sklearn.cluster.MiniBatchKMeans.html

| Movies | Genres | | | Keywords | | | Ratings | | |
|---|---|---|---|---|---|---|---|---|---|
| | D1 | D2 | D3 | D1 | D2 | D3 | D1 | D2 | D3 |
| Space | 0.301 | 0.358 | 0.354 | 0.185 | 0.198 | 0.201 | 0.463 | 0.475 | 0.486 |
| Single directions | **0.436** | 0.463 | 0.492 | 0.23 | **0.233** | **0.224** | 0.466 | 0.499 | 0.498 |
| Clusters | 0.431 | **0.513** | **0.506** | 0.215 | 0.22 | 0.219 | **0.504** | **0.507** | **0.513** |
| PPMI | 0.429 | 0.443 | 0.483 | **0.243** | 0.224 | 0.224 | 0.47 | 0.453 | 0.453 |
| Topic | 0.415 | 0.472 | 0.455 | 0.189 | 0.05 | 0.075 | 0.473 | 0.243 | 0.38 |

| Newsgroups | | | | Sentiment | | | Reuters | | |
|---|---|---|---|---|---|---|---|---|---|
| | D1 | D2 | D3 | D1 | D2 | D3 | D1 | D2 | D3 |
| Rep | 0.251 | 0.366 | 0.356 | 0.705 | 0.77 | 0.773 | 0.328 | 0.413 | 0.501 |
| Single dir | 0.418 | **0.49** | **0.537** | 0.784 | 0.814 | **0.821** | **0.678** | **0.706** | 0.72 |
| Cluster | 0.394 | 0.433 | 0.513 | 0.735 | **0.844** | 0.813 | 0.456 | 0.569 | 0.583 |
| PPMI | 0.33 | 0.407 | 0.444 | 0.7 | 0.719 | 0.73 | 0.616 | 0.699 | **0.723** |
| Topic | **0.431** | 0.423 | 0.444 | **0.79** | 0.791 | 0.811 | 0.411 | 0.527 | 0.536 |

| Placetypes | Foursquare | | | OpenCYC | | | Geonames | | |
|---|---|---|---|---|---|---|---|---|---|
| | D1 | D2 | D3 | D1 | D2 | D3 | D1 | D2 | D3 |
| Rep | 0.438 | 0.478 | 0.454 | 0.383 | 0.397 | 0.396 | 0.349 | 0.34 | 0.367 |
| Single dir | **0.541** | 0.498 | **0.531** | 0.404 | **0.428** | 0.39 | **0.444** | **0.533** | **0.473** |
| Cluster | 0.462 | 0.507 | 0.496 | **0.413** | 0.42 | **0.429** | 0.444 | 0.458 | 0.47 |
| PPMI | 0.473 | **0.512** | 0.491 | 0.371 | 0.351 | 0.352 | 0.361 | 0.301 | 0.242 |
| Topic | 0.488 | 0.433 | 0.526 | 0.365 | 0.271 | 0.313 | 0.365 | 0.3 | 0.219 |

**Table 3.2: summary of all results**

### 3.3.4 Baseline Representations

To begin, we show in Table 3.3 all variations of the baseline representations used directly as input to Decision Trees and SVM's. These examples that do not apply our methodology, serve as a reference point for what is possible using standard linear models without the need for interpretability. There is a big performance drop when going from depth three trees to depth one trees. These kind of performance drops are expected for these representations, as they do not have dimensions that correspond to key semantics, so it is unlikely that a smaller tree will be able to use the available dimensions to classify well. In this full table we include the precision and recall scores for clarity, mainly to explain why the high recall scores occur. This is because we balanced the weights as one of our hyper-parameters, and when the weight is balanced so that positive instances are weighted more heavily, the model prioritizes recall over precision. When this high recall score doesn't occur, that means that not balancing the weights performed better on the development data.

The size of the space is not as influential as the representation type in these results for the Decision Trees. For this reason we show only the best performing representation of each type in the main results table for this section. Although Linear SVM's perform the best on these representations without the need for interpretability, future results will be for low-depth decision trees in-order to easily distinguish the degree to which key semantics correspond to dimensions in the representations.

| Newsgroups | D1 | | | | D2 | | | | D3 | | | | DN | | | | SVM | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec |
| PCA 200 | 0.701 | 0.251 | 0.148 | 0.811 | 0.843 | 0.366 | 0.245 | 0.719 | 0.956 | 0.355 | 0.54 | 0.265 | 0.946 | 0.44 | 0.45 | 0.432 | 0.969 | 0.612 | 0.746 | 0.519 |
| PCA 100 | 0.698 | 0.247 | 0.146 | 0.813 | 0.835 | 0.362 | 0.241 | 0.731 | 0.957 | 0.356 | 0.576 | 0.257 | 0.948 | 0.451 | 0.465 | 0.438 | 0.969 | 0.586 | 0.768 | 0.474 |
| PCA 50 | 0.68 | 0.24 | 0.141 | 0.829 | 0.834 | 0.355 | 0.234 | 0.735 | 0.957 | 0.329 | 0.472 | 0.253 | 0.947 | 0.45 | 0.462 | 0.438 | 0.966 | 0.52 | 0.745 | 0.399 |
| AWV 200 | 0.687 | 0.217 | 0.126 | 0.781 | 0.758 | 0.256 | 0.156 | 0.718 | 0.764 | 0.26 | 0.157 | 0.751 | 0.937 | 0.339 | 0.352 | 0.328 | 0.961 | 0.468 | 0.641 | 0.369 |
| AWV 100 | 0.677 | 0.21 | 0.122 | 0.775 | 0.78 | 0.275 | 0.173 | 0.683 | 0.746 | 0.25 | 0.149 | 0.769 | 0.934 | 0.324 | 0.332 | 0.317 | 0.865 | 0.4 | 0.265 | 0.812 |
| AWV 50 | 0.696 | 0.219 | 0.127 | 0.772 | 0.777 | 0.272 | 0.168 | 0.71 | 0.743 | 0.25 | 0.149 | 0.786 | 0.935 | 0.325 | 0.335 | 0.316 | 0.842 | 0.362 | 0.233 | 0.819 |
| MDS 200 | 0.581 | 0.184 | 0.103 | **0.837** | 0.742 | 0.262 | 0.16 | 0.729 | 0.719 | 0.236 | 0.139 | 0.785 | 0.935 | 0.327 | 0.332 | 0.323 | 0.965 | 0.501 | **0.802** | 0.364 |
| MDS 100 | 0.586 | 0.187 | 0.105 | 0.833 | 0.754 | 0.261 | 0.159 | 0.727 | 0.705 | 0.236 | 0.138 | **0.808** | 0.935 | 0.33 | 0.338 | 0.321 | 0.878 | 0.439 | 0.308 | 0.765 |
| MDS 50 | 0.593 | 0.153 | 0.087 | 0.647 | 0.716 | 0.25 | 0.15 | **0.756** | 0.736 | 0.243 | 0.144 | 0.774 | 0.935 | 0.324 | 0.335 | 0.313 | 0.854 | 0.394 | 0.259 | 0.821 |
| D2V 200 | 0.682 | 0.205 | 0.119 | 0.746 | 0.802 | 0.268 | 0.169 | 0.646 | 0.77 | 0.269 | 0.164 | 0.75 | 0.94 | 0.366 | 0.389 | 0.346 | 0.961 | 0.468 | 0.641 | 0.369 |
| D2V 100 | 0.682 | 0.208 | 0.12 | 0.762 | 0.792 | 0.268 | 0.168 | 0.662 | 0.786 | 0.268 | 0.164 | 0.727 | 0.94 | 0.376 | 0.392 | 0.361 | **0.971** | **0.628** | 0.761 | 0.535 |
| D2V 50 | 0.683 | 0.207 | 0.12 | 0.764 | 0.809 | 0.294 | 0.187 | 0.694 | 0.782 | 0.28 | 0.172 | 0.761 | 0.943 | 0.394 | 0.415 | 0.376 | 0.97 | 0.601 | 0.758 | 0.497 |
| PPMI | **0.948** | 0.33 | **0.532** | 0.239 | 0.947 | 0.407 | 0.511 | 0.338 | 0.944 | **0.444** | 0.506 | 0.396 | **0.951** | **0.494** | **0.496** | **0.492** | 0.962 | 0.613 | 0.627 | 0.599 |
| Topic | 0.852 | **0.431** | 0.304 | 0.743 | **0.96** | **0.423** | **0.604** | 0.326 | **0.961** | 0.444 | **0.606** | 0.35 | 0.944 | 0.432 | 0.434 | 0.429 | 0.879 | 0.46 | 0.318 | **0.835** |

**Table 3.3: Full results for the newsgroups.**

**Table 3.4: Results for all other domains for the representations.**

**Reuters**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.847 | 0.328 | 0.917 | 0.413 | 0.978 | 0.501 | 0.978 | 0.565 | 0.989 | 0.761 |
| AWV | 0.782 | 0.252 | 0.971 | 0.328 | 0.974 | 0.417 | 0.973 | 0.495 | 0.987 | 0.719 |
| MDS | 0.791 | 0.263 | 0.9 | 0.357 | 0.979 | 0.489 | 0.976 | 0.522 | 0.988 | 0.67 |
| D2V | 0.818 | 0.268 | 0.867 | 0.298 | 0.974 | 0.445 | 0.971 | 0.482 | 0.986 | 0.724 |
| PPMI | **0.975** | **0.616** | **0.978** | **0.699** | **0.98** | **0.723** | **0.984** | **0.746** | **0.99** | **0.8** |
| Topic | 0.92 | 0.411 | 0.977 | 0.527 | 0.977 | 0.536 | 0.977 | 0.56 | 0.95 | 0.513 |

**Sentiment**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.745 | 0.705 | 0.755 | 0.77 | 0.778 | 0.773 | **0.781** | **0.779** | **0.891** | **0.893** |
| AWV | 0.642 | 0.652 | 0.643 | 0.694 | 0.695 | 0.717 | 0.66 | 0.663 | 0.827 | 0.829 |
| D2V | 0.642 | 0.664 | 0.66 | 0.707 | 0.702 | 0.7 | 0.711 | 0.708 | 0.878 | 0.878 |
| PPMI | 0.616 | 0.7 | 0.655 | 0.719 | 0.675 | 0.73 | 0.712 | 0.71 | 0.887 | 0.888 |
| Topic | **0.793** | **0.79** | **0.794** | **0.791** | **0.81** | **0.811** | 0.733 | 0.73 | 0.815 | 0.822 |

**Placetypes OpenCYC**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.586 | 0.346 | 0.708 | 0.343 | 0.695 | 0.342 | 0.832 | 0.309 | 0.847 | 0.474 |
| AWV | 0.625 | **0.383** | 0.651 | 0.376 | 0.728 | **0.396** | **0.844** | **0.362** | 0.85 | 0.466 |
| MDS | 0.624 | 0.364 | 0.7 | **0.397** | 0.731 | 0.374 | 0.843 | 0.305 | 0.861 | **0.476** |
| PPMI | **0.728** | 0.371 | 0.75 | 0.351 | 0.739 | 0.352 | 0.843 | 0.323 | **0.9** | 0.366 |
| Topic | 0.708 | 0.365 | **0.87** | 0.271 | **0.87** | 0.313 | 0.831 | 0.313 | 0.808 | 0.407 |

**Movies Genres**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.722 | 0.301 | 0.755 | 0.339 | 0.717 | 0.321 | 0.884 | 0.372 | **0.925** | 0.518 |
| AWV | 0.679 | 0.29 | 0.774 | 0.321 | 0.756 | 0.343 | 0.873 | 0.312 | 0.922 | 0.496 |
| MDS | 0.679 | 0.298 | 0.79 | 0.358 | 0.773 | 0.354 | 0.887 | 0.385 | 0.875 | **0.532** |
| PPMI | **0.852** | **0.429** | **0.91** | 0.443 | **0.912** | **0.483** | 0.882 | **0.416** | 0.923 | 0.526 |
| Topic | 0.767 | 0.415 | 0.905 | **0.472** | 0.912 | 0.455 | **0.889** | 0.415 | 0.843 | 0.491 |

**Placetypes Foursquare**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.731 | 0.342 | 0.823 | 0.393 | 0.86 | 0.388 | 0.887 | 0.398 | 0.896 | 0.568 |
| AWV | 0.767 | 0.401 | 0.828 | 0.478 | 0.85 | 0.452 | 0.905 | **0.505** | 0.923 | **0.622** |
| MDS | **0.915** | 0.438 | 0.804 | 0.427 | 0.86 | 0.454 | 0.893 | 0.462 | 0.932 | 0.619 |
| PPMI | 0.889 | 0.473 | 0.915 | **0.512** | 0.904 | 0.491 | 0.881 | 0.31 | **0.938** | 0.567 |
| Topic | 0.864 | **0.488** | **0.916** | 0.433 | **0.917** | **0.526** | **0.907** | 0.464 | 0.916 | 0.569 |

**Movies Keywords**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.647 | 0.185 | 0.644 | 0.193 | 0.677 | 0.199 | 0.846 | 0.161 | 0.787 | 0.272 |
| AWV | 0.5 | 0.16 | 0.641 | 0.174 | 0.595 | 0.174 | 0.853 | 0.141 | 0.717 | 0.23 |
| MDS | 0.633 | 0.179 | 0.69 | 0.201 | 0.674 | 0.201 | 0.84 | 0.163 | 0.788 | **0.28** |
| PPMI | **0.818** | **0.243** | 0.745 | **0.224** | 0.739 | **0.224** | 0.847 | **0.17** | **0.921** | 0.217 |
| Topic | 0.629 | 0.189 | **0.932** | 0.05 | **0.93** | 0.075 | **0.857** | 0.152 | 0.678 | 0.21 |

**Placetypes Geonames**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | 0.502 | 0.301 | 0.69 | 0.305 | 0.68 | 0.295 | 0.821 | 0.243 | 0.844 | 0.401 |
| AWV | 0.657 | 0.326 | 0.755 | 0.323 | 0.842 | **0.367** | 0.813 | 0.332 | 0.865 | **0.514** |
| MDS | 0.626 | 0.349 | 0.695 | **0.34** | 0.796 | 0.272 | **0.845** | 0.295 | 0.638 | 0.397 |
| PPMI | **0.808** | 0.361 | 0.732 | 0.301 | 0.76 | 0.242 | 0.83 | 0.283 | **0.894** | 0.312 |
| Topic | 0.771 | **0.365** | **0.863** | 0.3 | **0.85** | 0.219 | 0.828 | **0.348** | 0.819 | 0.349 |

**Movies Ratings**

| | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | DN ACC | DN F1 | SVM ACC | SVM F1 |
|---|---|---|---|---|---|---|---|---|---|---|
| PCA | **0.65** | 0.463 | 0.681 | **0.475** | 0.684 | **0.486** | 0.744 | 0.408 | 0.771 | 0.58 |
| AWV | 0.601 | 0.423 | 0.618 | 0.433 | 0.596 | 0.448 | 0.736 | 0.372 | 0.73 | 0.532 |
| MDS | 0.592 | 0.437 | 0.635 | 0.449 | 0.631 | 0.452 | **0.752** | **0.412** | 0.773 | **0.589** |
| PPMI | 0.583 | 0.47 | 0.635 | 0.453 | 0.605 | 0.453 | 0.73 | 0.384 | **0.825** | 0.536 |
| Topic | 0.575 | **0.473** | **0.789** | 0.243 | **0.789** | 0.38 | 0.739 | 0.375 | 0.704 | 0.501 |

### 3.3.5   Semantic Spaces

In this section, we explain how we obtained four different Semantic Spaces.

### 3.3.6   Word Directions

For all trees we use grid search to find the best values for the criterion, either the gini score or the information entropy score, the maximum amount of features between [None, 'auto', 'log2'], and additionally, we include whether or not to balance the classes in the grid search.

These single directions typically overfit.

| Newsgroups | D1 | | | | D2 | | | | D3 | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec |
| PCA 200 | 0.955 | 0.348 | 0.521 | 0.261 | 0.959 | 0.424 | 0.678 | 0.309 | 0.96 | 0.454 | 0.674 | 0.343 |
| PCA 100 | 0.957 | 0.382 | 0.491 | 0.313 | 0.961 | 0.474 | 0.679 | 0.364 | **0.963** | 0.512 | 0.694 | 0.406 |
| PCA 50 | 0.957 | 0.373 | 0.417 | 0.337 | **0.963** | 0.478 | 0.621 | 0.388 | 0.963 | 0.506 | 0.7 | 0.396 |
| AWV 200 | 0.832 | 0.35 | 0.226 | 0.777 | 0.957 | 0.383 | 0.517 | 0.305 | 0.958 | 0.445 | 0.598 | 0.354 |
| AWV 100 | 0.83 | 0.343 | 0.219 | 0.785 | 0.823 | 0.36 | 0.233 | **0.792** | 0.956 | 0.387 | 0.563 | 0.295 |
| AWV 50 | 0.807 | 0.341 | 0.215 | 0.816 | 0.833 | 0.361 | 0.236 | 0.762 | 0.954 | 0.392 | 0.511 | 0.318 |
| MDS 200 | **0.959** | **0.418** | **0.543** | 0.339 | 0.962 | 0.465 | 0.669 | 0.357 | 0.962 | 0.493 | **0.707** | 0.379 |
| MDS 100 | 0.857 | 0.365 | 0.244 | 0.725 | 0.959 | 0.428 | 0.624 | 0.326 | 0.96 | 0.453 | 0.644 | 0.349 |
| MDS 50 | 0.821 | 0.324 | 0.206 | 0.762 | 0.842 | 0.386 | 0.258 | 0.77 | 0.957 | 0.398 | 0.596 | 0.299 |
| D2V 200 | 0.831 | 0.343 | 0.22 | 0.784 | 0.96 | 0.47 | **0.683** | 0.358 | 0.962 | 0.494 | 0.69 | 0.385 |
| D2V 100 | 0.844 | 0.374 | 0.243 | 0.803 | 0.961 | **0.49** | 0.642 | 0.396 | 0.962 | 0.517 | 0.67 | 0.421 |
| D2V 50 | 0.845 | 0.388 | 0.252 | **0.844** | 0.962 | 0.488 | 0.639 | 0.395 | 0.963 | **0.537** | 0.673 | **0.446** |

**Table 3.5: Newsgroups single dirs**

| Newsgroups | D1 | | | | D2 | | | | D3 | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec | ACC | F1 | Prec | Rec |
| PCA 200 | 0.955 | 0.348 | 0.521 | *0.261* | 0.959 | 0.424 | 0.678 | 0.309 | 0.96 | 0.454 | 0.674 | 0.343 |
| PCA 100 | 0.957 | 0.382 | 0.491 | 0.313 | 0.961 | 0.474 | 0.679 | 0.364 | **0.963** | 0.512 | 0.694 | 0.406 |
| PCA 50 | 0.957 | 0.373 | 0.417 | 0.337 | **0.963** | 0.478 | 0.621 | 0.388 | 0.963 | 0.506 | 0.7 | 0.396 |
| AWV 200 | 0.832 | 0.35 | 0.226 | 0.777 | 0.957 | 0.383 | 0.517 | *0.305* | 0.958 | 0.445 | 0.598 | 0.354 |
| AWV 100 | 0.83 | 0.343 | 0.219 | 0.785 | *0.823* | *0.36* | *0.233* | **0.792** | 0.956 | *0.387* | 0.563 | *0.295* |
| AWV 50 | *0.807* | 0.341 | 0.215 | 0.816 | 0.833 | 0.361 | 0.236 | 0.762 | *0.954* | 0.392 | *0.511* | 0.318 |
| MDS 200 | **0.959** | **0.418** | **0.543** | 0.339 | 0.962 | 0.465 | 0.669 | 0.357 | 0.962 | 0.493 | **0.707** | 0.379 |
| MDS 100 | 0.857 | 0.365 | 0.244 | 0.725 | 0.959 | 0.428 | 0.624 | 0.326 | 0.96 | 0.453 | 0.644 | 0.349 |
| MDS 50 | 0.821 | *0.324* | *0.206* | 0.762 | 0.842 | 0.386 | 0.258 | 0.77 | 0.957 | 0.398 | 0.596 | 0.299 |
| D2V 200 | 0.831 | 0.343 | 0.22 | 0.784 | 0.96 | 0.47 | **0.683** | 0.358 | 0.962 | 0.494 | 0.69 | 0.385 |
| D2V 100 | 0.844 | 0.374 | 0.243 | 0.803 | 0.961 | **0.49** | 0.642 | 0.396 | 0.962 | 0.517 | 0.67 | 0.421 |
| D2V 50 | 0.845 | 0.388 | 0.252 | **0.844** | 0.962 | 0.488 | 0.639 | 0.395 | 0.963 | **0.537** | 0.673 | **0.446** |

| Reuters | D1 | | D2 | | D3 | |
|---|---|---|---|---|---|---|
| | ACC | F1 | ACC | F1 | ACC | F1 |
| PCA | 0.976 | 0.658 | *0.979* | 0.679 | *0.977* | *0.467* |
| AWV | *0.975* | 0.598 | 0.979 | *0.656* | 0.98 | 0.66 |
| MDS | 0.975 | **0.678** | **0.98** | **0.706** | **0.982** | **0.72** |
| D2V | **0.977** | *0.583* | 0.979 | 0.664 | 0.98 | 0.632 |

| Sentiment | D1 | | D2 | | D3 | |
|---|---|---|---|---|---|---|
| | ACC | F1 | ACC | F1 | ACC | F1 |
| PCA | 0.739 | 0.759 | **0.797** | **0.814** | 0.802 | 0.805 |
| AWV | *0.7* | *0.699* | *0.711* | *0.736* | *0.723* | *0.735* |
| D2V | **0.776** | **0.784** | 0.782 | 0.801 | **0.822** | **0.821** |

| Placetypes | D1 | | D2 | | D3 | |
|---|---|---|---|---|---|---|
| | ACC | F1 | ACC | F1 | ACC | F1 |
| OpenCYC | | | | | | |
| PCA | *0.632* | *0.371* | *0.704* | *0.381* | *0.735* | 0.365 |

| Movies | D1 | | D2 | | D3 | |
|---|---|---|---|---|---|---|
| | ACC | F1 | ACC | F1 | ACC | F1 |
| Genres | | | | | | |
| PCA | 0.824 | 0.441 | *0.82* | *0.412* | 0.913 | 0.463 |

### 3.3.7 Clustered Directions

**Table 3.7: All clustering results**

**Newsgroups**

| Newsgroups | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| K-means 200 | **0.852** | **0.394** | **0.261** | 0.795 | **0.958** | **0.433** | **0.58** | **0.963** | **0.513** | **0.704** | 0.345 | 0.403 |
| K-means 100 | 0.842 | 0.388 | 0.257 | 0.791 | 0.958 | 0.366 | 0.516 | 0.962 | 0.5 | 0.635 | 0.284 | **0.412** |
| K-means 50 | 0.834 | 0.381 | 0.248 | **0.819** | 0.815 | 0.336 | 0.212 | 0.961 | 0.485 | 0.612 | **0.81** | 0.402 |
| Derrac 200 | 0.803 | 0.313 | 0.202 | 0.693 | 0.797 | 0.306 | 0.191 | 0.958 | 0.409 | 0.605 | 0.781 | 0.309 |
| Derrac 100 | 0.792 | 0.305 | 0.197 | 0.667 | 0.791 | 0.287 | 0.179 | 0.957 | 0.374 | 0.56 | 0.721 | 0.281 |
| Derrac 50 | 0.769 | 0.26 | 0.162 | 0.661 | 0.768 | 0.237 | 0.143 | 0.955 | 0.315 | 0.47 | 0.693 | 0.237 |

**Reuters / Sentiment**

| | Reuters D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | Sentiment D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| K-means | **0.875** | **0.338** | **0.975** | **0.54** | 0.973 | **0.58** | 0.623 | 0.674 | **0.837** | **0.844** | 0.658 | 0.707 |
| Derrac | 0.797 | 0.291 | 0.973 | 0.402 | **0.974** | 0.485 | **0.712** | **0.735** | 0.802 | 0.82 | **0.803** | **0.813** |

**Placetypes (OpenCYC) / Movies (Genres)**

| | OpenCYC D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | Genres D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| K-means | **0.641** | **0.413** | **0.735** | **0.405** | **0.813** | **0.431** | **0.913** | **0.511** | **0.913** | **0.513** | **0.913** | **0.506** |
| Derrac | 0.605 | 0.39 | 0.672 | 0.392 | 0.759 | 0.341 | 0.805 | 0.425 | 0.789 | 0.431 | 0.911 | 0.432 |

**Foursquare (Keywords) / Geonames (Ratings)**

| | Keywords D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | Ratings D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| K-means | **0.913** | **0.462** | **0.911** | **0.5** | **0.891** | **0.511** | 0.667 | 0.208 | 0.648 | 0.202 | 0.678 | 0.213 |
| Derrac | 0.768 | 0.392 | 0.835 | 0.445 | 0.805 | 0.425 | **0.726** | **0.215** | **0.745** | **0.22** | **0.707** | **0.219** |

**Geonames / Ratings**

| | Geonames D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 | Ratings D1 ACC | D1 F1 | D2 ACC | D2 F1 | D3 ACC | D3 F1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| K-means | **0.772** | 0.43 | **0.774** | 0.407 | **0.819** | **0.472** | **0.671** | **0.504** | 0.638 | **0.507** | **0.686** | **0.513** |
| Derrac | 0.678 | **0.449** | 0.74 | **0.411** | 0.807 | 0.415 | 0.651 | 0.445 | **0.669** | 0.463 | 0.627 | 0.479 |

## 3.4 Qualitative Results

Make reference to the qualitative results found in the previous work here.

### 3.4.1 Examining the differences between directions

Investigating three potential hypothesis: 1. The ranks are more accurate, so the key directions are better represented that would contribute 2. The spaces/score-types contain unique directions that contribute to the tree directly 3. The spaces/score-types influence the rankings so that they are better represented, but are not directly used in the tree

First, we look at the best scoring directions. Then, we look at the unique directions for each space-type and score-type. The section is then followed by conclusions, and we begin to look into the clusters.

### 3.4.2 The best-performing directions for each space type

What are the domains that best convey the similarities and differences between different domains?

1. Find domains that act differently (perhaps one domain where a space-type that is not usually scoring high is scoring high, big differences in F1) 2. Get interesting directions from those domains

| Movies (50 MDS NDCG) | Sentiment (100 D2V NDCG) | Newsgroups (50 D2V NDCG) | Place-types (50 PCA Kappa) | Reuters (200 MDS NDCG) |
|---|---|---|---|---|
| horror (scares, scary) | glenda (glen, matthau) | karabag (iranian, turkiye) | blackcountry (listed, westmidlands) | franklin (fund, mthly) |
| hilarious (funniest, hilarity) | scarlett (gable, dalton) | leftover (flaming, vancouver) | ears (stare, adorable) | quarterly (shearson, basis) |
| bollywood (hindi, india) | giallo (argento, fulci) | wk (5173552178, 18084tmibmclmsuedu) | spagna (espanha, colores) | feb (28, splits) |
| laughs (funnier, funniest) | bourne (damon, cusack) | 1069 (mlud, wibbled) | oldfashioned (winery, antiques) | 22 (booked, hong) |
| jokes (gags, laughs) | piper (omen, knightley) | providence (norris, ahl) | gardening (greenhouse, petals) | april (monthly, average) |
| comedies (comedic, laughs) | casper (dolph, damme) | celestial (interplanetary, bible) | pagoda (hindu, carved) | sets (principally, precious) |
| hindi (bollywood, india) | norris (chuck, rangers) | mlud (wibbled, 1069) | artificial (saturation, cs4) | 16 (creditor, trillion) |
| war (military, army) | holmes (sherlock, rathbone) | endif (olwm, ciphertext) | inner (curved, rooftops) | 1st (qtr, pennsylvania) |
| western (outlaw, unforgiven) | rourke (mickey, walken) | gd3004 (35894, intergraph) | celebrate (festive, celebrity) | 26 (approve, inadequate) |
| romantic (romance, chemistry) | ustinov (warden, cassavetes) | rtfmmitedu (newsanswers, ieee) | vietnamese (ethnic, hindu) | 23 (offsetting, weekly) |
| songs (song, tunes) | scooby (doo, garfield) | eng (padres, makefile) | cn (elevated, amtrak) | prior (recapitalization, payment) |
| sci (science, outer) | doo (scooby, garfield) | pizza (bait, wiretap) | mannequin (bags, jewelry) | avg (shrs, shr) |
| funniest (hilarious, funnier) | heston (charlton, palance) | porsche (nanao, mercedes) | falcon (r, 22) | june (july, venice) |
| noir (noirs, bogart) | homer (pacino, macy) | gebcadredslpittedu (n3jxp, skepticism) | jewish (monuments, cobblestone) | march (31, day) |
| documentary (documentaries, footage) | welles (orson, kane) | scsi2 (scsi, cooling) | canon60d (kitlens, 600d) | regular (diesel, petrol) |
| animation (animated, animators) | frost (snowman, damme) | playback (quicktime, xmotif) | reflective (curved, cropped) | 4th (qtr, fourth) |
| adults (adult, children) | streisand (bridget, salman) | 35894 (gd3004, medin) | mason (edward, will) | 27 (chemlawn, theyre) |
| creepy (spooky, scary) | davies (rhys, marion) | diesel (volvo, shotguns) | aerialview (manmade, largest) | 14 (borrowing, borrowings) |
| gay (gays, homosexuality) | cinderella (fairy, stepmother) | evolutionary (shifting, hulk) | shelf (rack, boxes) | 11 (chapter, ranged) |
| workout (intermediate, instruction) | boll (uwe, belushi) | techniciandr (obp, 144k) | monroe (raleigh, jefferson) | may (probably, however) |
| thriller (thrillers, suspense) | rochester (eyre, dalton) | 8177 (obp, 144k) | litter (fujichrome, e6) | 38 (33, strong) |
| funnier (laughs, funniest) | edie (soprano, vertigo) | shaw (medicine, ottoman) | streetlights (streetlamp, headlights) | m1 (m2, m3) |
| suspense (suspenseful, thrillers) | scarecrow (zombies, reese) | scorer (gilmour, lindros) | carlzeiss (f2, voigtlander) | dlr (writedown, debt) |
| arts (hong, chan) | kramer (streep, meryl) | xwd (xloadimage, openwindows) | manmade (aerialview, below) | five (years, jones) |
| christianity (religious, religion) | marty (amitabh, goldie) | ee (275, xloadimage) | demolished (neglected, rundown) | bushels (soybeans, ccc) |
| musical (singing, sing) | columbo (falk, garfield) | com2 (com1, v32bis) | wald (berge, wildflower) | revs (net, 3for2) |
| gore (gory, blood) | kidman (nicole, jude) | examiner (corpses, brass) | arquitetura (exposition, cidade) | 29 (175, include) |
| animated (animation, cartoon) | juliet (romeo, troma) | migraine (ama, placebo) | greyscale (highcontrast, monochromatic) | acquisition (make, usairs) |
| gags (jokes, slapstick) | garland (judy, lily) | parliament (parliamentary, armored) | alameda (monday, marin) | payable (div, close) |

**Table 3.8: Table**

### 3.4.3   How Domain Directions Differ

For the single directions, arrange them by score where the highest scoring directions are at the top. For the clusters, there is no convenient way to organize them without bias, so clusters that are interesting are selected.

**Score Types**

There are unique directions for each different space type, each suitable to different tasks. NDCG was selected as the best score-type for Sentiment, Newsgroups, Reuters, Movies Genres, Movies Keywords in depth-3 Decision Trees. Place-types foursquare used F1-score, but the classes are very unbalanced and there are few documents.

| NDCG | F1 | Accuracy | Kappa | Common |
|---|---|---|---|---|
| gay (homosexuality, sexuality) | company (sell, pay) | kennedy (republic, elected) | definately (alot, awesome) | horror (scares, scares) |
| arts (hong, chan) | street (city, york) | bags (listened, salvation) | guns (gun, shoot) | laughs (funnier, funnier) |
| sports (win, players) | red (numerous, fashion) | summers (verge, medieval) | flawless (perfection, brilliantly) | jokes (gags, gags) |
| apes (remembered, planet) | project (creating, spent) | revolve (sincerely, historian) | mail (reviewed, rated) | comedies (comedic, comedic) |
| german (germans, europe) | mark (favor, pull) | locale (foster, sharply) | garbage (crap, horrible) | sci (scifi, alien) |
| satire (parody, parodies) | lady (actress, lovely) | cooler (downward, reports) | featurette (featurettes, extras) | funniest (hilarious, hilarious) |
| band (rock, vocals) | fire (ground, force) | spades (ralph, medieval) | complaint (extra, added) | creepy (spooky, spooky) |
| crude (offensive, offended) | post (essentially, purpose) | filmography (ralph, experiments) | mission (enemy, saving) | thriller (thrillers, thrillers) |
| dancing (dance, dances) | heads (large, throw) | quentin (downward, anime) | ruin (wondering, heck) | funnier (laughs, laughs) |
| restored (print, remastered) | water (land, large) | employers (finishes, downward) | wars (forces, enemy) | suspense (suspenseful, suspenseful) |
| drugs (drug, abuse) | road (drive, trip) | formal (victory, kennedy) | prefer (compare, added) | gore (gory, gory) |
| church (religious, jesus) | brother (son, dad) | tube (esta, muscle) | heroes (packed, hero) | gags (jokes, jokes) |
| sexuality (sexual, sexually) | party (decide, hot) | woefully (restless, knockout) | necessarily (offer, draw) | science (sci, sci) |
| sexually (sexual, sexuality) | badly (awful, poorly) | scientists (hilarity, locale) | portray (portrayed, portraying) | gory (gore, gore) |
| england (british, english) | limited (aspect, unlike) | overboard (civilized, cinderella) | critic (reviewed, net) | government (political, political) |
| ocean (sea, boat) | impression (instance, reasons) | rumors (homosexuality, characteristics) | reviewed (rated, mail) | suspenseful (suspense, suspense) |
| marry (married, marriage) | trip (journey, road) | salvation (bags, cooler) | saving (carry, forced) | frightening (terrifying, terrifying) |
| campy (cult, cheesy) | michael (producers, david) | actively (assassination, overcoming) | technical (digital, presentation) | military (army, army) |
| christian (religious, jesus) | memory (forgotten, memories) | stretching (victory, hideous) | statement (exist, critical) | slapstick (gags, gags) |
| melodrama (dramatic, tragedy) | james (robert, michael) | downward (cooler, crawling) | shocked (hate, warning) | scary (scare, scare) |
| sing (singing, sings) | thin (barely, flat) | rocked (staple, demented) | flying (air, force) | blu (unanswered, ray) |
| sentimental (touching, sappy) | pre (popular, include) | affectionate (esta, muscle) | danger (dangerous, edge) | internetreviews (rhodes, rhodes) |
| depressing (bleak, suffering) | faces (constant, unlike) | protest (protective, assassination) | | cgi (computer, computer) |
| evidence (investigation, accused) | values (exception, wise) | confined (cooler, downward) | | email (web, web) |
| adorable (cute, sweet) | unusual (odd, seemingly) | inhabit (quentin, drawback) | | thrilling (thrill, exciting) |
| episodes (episode, television) | lovers (lover, lovely) | latin (communities, mount) | | web (email, email) |
| teenager (teen, teenage) | frame (image, effect) | reception (como, finishes) | | horror (scares, scares) |
| magical (fantasy, lovely) | mans (ultimate, sees) | uptight (suspensful, stalked) | | laughs (funnier, funnier) |
| health (medical, suffering) | efforts (generally, nonetheless) | brink (inexplicable, freddy) | | suspense (suspenseful, suspenseful) |

**Table 3.9: Different score types**

**Comparing Space Types**

We begin by selecting the space that performed well on the genres task for the movies, with the understanding that genres as a key natural classification task will likely make use of good directions that correspond to domain knowledge. After selecting this space, we choose similarly sized spaces from the other space-types, in this case we selected the 200 dimensional MDS space as it performed the best and from there, we selected the 200 dimensional PCA space and AWV space. We also use the same score-type and frequency cut-off as the best performing space-type. In this case, the best performing type for the PCA space was 20000 frequency cutoff and NDCG, and we are comparing to 10000 frequency cutoff. This means that we are sometimes using a slightly worse performing space-type than the one we used as our final results, and that the original space has a performance advantage, but we have chosen to do so to make the results more consistent and specific. We approach these qualitative experiments with the following idea: spaces that perform better on natural domain tasks using decision trees contain unique natural directions that other spaces do not have.

The commonalities between spaces are much more prevalent than the differences, with natural concepts of the domain being represented in all of the different space types. However, different spaces do perform better than others on natural domain tasks. In this section, we investigate why this occurs and the differences between spaces built using a standard frequency-based approach, word-vectors and doc2vec, which uses a combination of contextual information and word vectors.

**Comparing MDS, AWV and PCA in the Movies domain**

| MDS | AWV | PCA | Common |
|---|---|---|---|
| berardinelli (employers, distributor) | billy (thrown, dirty) | amount (leaving, pick) | noir (fatale, femme) |
| crawford (joan, davis) | brother (brothers, boys) | fails (fit, pick) | gay (homosexual, homosexuality) |
| hitchcocks (hitchcock, alfred) | fonda (henry, jane) | pick (fails, fit) | prison (jail, prisoners) |
| warners (warner, bros) | building (built, climax) | stands (fails, cover) | arts (rec, robomod) |
| nuclear (weapons, soviet) | train (tracks, thrown) | surprisingly (offer, fit) | allens (woody, allen) |
| joan (crawford, barbara) | slaves (slavery, excuse) | copyright (email, compuserve) | jokes (laughs, joke) |
| kidnapped (kidnapping, torture) | | length (reflect, expressed) | animation (animated, cartoon) |
| hop (hip, rap) | | profanity (reflect, producers) | sherlock (holmes, detective) |
| kung (martial, jackie) | | compuserve (copyright, internetreviews) | western (westerns, wayne) |
| ballet (dancers, dancer) | | talents (admit, agree) | songs (song, lyrics) |
| gambling (vegas, las) | | admit (agree, talents) | comedies (comedic, laughs) |
| alcoholic (drunk, alcoholism) | | developed (introduced, sounds) | workout (exercise, challenging) |
| waves (surfing, wave) | | intended (bother, werent) | laughs (funnier, hilarious) |
| jaws (jurassic, godfather) | | constantly (putting, sounds) | drug (drugs, addict) |
| jungle (natives, island) | | tired (anymore, mediocre) | sci (science, fiction) |
| employers (berardinelli, distributor) | | produced (spoiler, surprising) | documentary (documentaries, interviews) |
| pot (weed, stoned) | | involving (believes, belief) | students (student, schools) |
| canadian (invasion, cheap) | | anymore (continue, tired) | thriller (thrillers, suspense) |
| murphy (eddie, comedian) | | leaving (fit, pick) | allen (woody, allens) |
| comics (comedian, comedians) | | makers (producers, aspects) | funniest (hilarious, laughing) |
| kidnapping (kidnapped, torture) | | introduced (developed, considered) | gags (jokes, slapstick) |
| subscribe (email, internetreviews) | | loses (climax, suffers) | adults (children, adult) |
| vegas (las, gambling) | | negative (positive, bother) | animated (animation, cartoon) |
| distributor (berardinelli, employers) | | expressed (reflect, opinions) | dancing (dance, dances) |
| wave (waves, surfing) | | mildly (mediocre, forgettable) | teen (teenage, teens) |
| rhodes (internetreviews, email) | | helped (putting, allowed) | soldiers (soldier, army) |
| hippie (pot, sixties) | | reflect (expressed, opinions) | indie (independent, festival) |
| weed (pot, stoned) | | opinions (reflect, expressed) | suspense (suspenseful, thriller) |
| caribbean (pirates, island) | | frequently (occasionally, consistently) | creepy (scary, eerie) |
| eddie (murphy, comedian) | | content (agree, proves) | italian (italy, spaghetti) |
| sixties (beatles, hippie) | | allowed (helped, werent) | jews (jewish, nazis) |
| ... 8 More | | suffers (lacks, loses) | ... 1480 more |

**Table 3.10: ok dude**

**Comparing PPMI representations to doc2vec**

### 3.4.4   What is the value of different score-types?

### 3.4.5   Producing Semantic Spaces

We use unsupervised representation learning methods, with the intention to obtain a representation that represents all salient features of the domain and can adapt to a variety of tasks.

For the semantic space, we compute the Positive Pointwise Mutual Information (See **??**) scores for the Bag-Of-Words, and use that as input to a variety of different off-the-shelf dimensionality reduction algorithms. We explain these in further detail in Section **??**.

### 3.4.6   Quantitative Results

From a domain, e.g. movie reviews, where each document is a collection of reviews for a movie, we preprocess the text such that it is converted to lower-case, and non-alphanumeric characters are removed. From here, we remove standard English stop words using the NLTK library [**?**]. We show an example of a review's original and converted formats in Figure **??**. From this preprocessed corpus, we obtain a Bag-Of-Words where we count the frequency of each term $BOW_w f$, see 2.1.1.

The difference between single directions and clusters is best highlighted when comparing their use in simple interpretable classifiers. In figure **??** we demonstrate this.

1. Negative directions (e.g. church for horror) 2. Non-contextualized, non-direct ways of classifying, versus clustering which finds salient properties which almost directly correspond to these natural tasks.

### 3.4.7   Interpretability Results

| D2V | MDS | Common |
|---|---|---|
| leftover *(pizza, brake)* | hi *(folks, everyone)* | chastity *(shameful, soon)* |
| wk *(5173552178, 18084tmibmclmsuedu)* | looking *(spend, rather)* | n3jxp *(gordon, gebcadredslpittedu)* |
| eng *(padres, makefile)* | need *(needs, means)* | skepticism *(gebcadredslpittedu, n3jxp)* |
| porsche *(nanao, 1280x1024)* | post *(summary, net)* | anyone *(knows, else)* |
| diesel *(cylinders, steam)* | find *(couldnt, look)* | gebcadredslpittedu *(soon, gordon)* |
| scorer *(gilmour, lindros)* | hello *(kind, thank)* | intellect *(soon, gordon)* |
| parliament *(caucasus, semifinals)* | david *(yet, man)* | please *(respond, reply)* |
| atm *(padres, inflatable)* | got *(mine, youve)* | thanks *(responses, advance)* |
| cryptology *(attendees, bait)* | go *(take, lets)* | email *(via, address)* |
| intake *(calcium, mellon)* | question *(answer, answered)* | know *(let, far)* |
| 433 *(366, 313)* | interested *(including, products)* | get *(wait, trying)* |
| ghetto *(warsaw, gaza)* | list *(mailing, send)* | think *(important, level)* |
| lens *(lenses, ankara)* | sorry *(guess, hear)* | good *(luck, bad)* |
| rushdie *(sinless, wiretaps)* | heard *(ever, anything)* | shafer *(dryden, nasa)* |
| immaculate *(porsche, alice)* | cheers *(kent, instead)* | bobbeviceicotekcom *(manhattan, beauchaine)* |
| keenan *(lindros, bosnian)* | say *(nothing, anything)* | dryden *(shafer, nasa)* |
| boxer *(jets, hawks)* | number *(call, numbers)* | im *(sure, working)* |
| linden *(mogilny, 176)* | mailing *(list, send)* | sank *(bronx, away)* |
| candida *(yeast, noring)* | call *(number, phone)* | banks *(soon, gordon)* |
| octopus *(web, 347)* | thank *(thanx, better)* | like *(sounds, looks)* |
| czech *(detectors, kuwait)* | read *(reading, group)* | shameful *(soon, gordon)* |
| survivor *(warsaw, croats)* | phone *(company, number)* | could *(away, bobbeviceicotekcom)* |
| 5173552178 *(circumference, wk)* | mail *(send, list)* | would *(appreciate, wouldnt)* |
| 18084tmibmclmsuedu *(circumference, wk)* | doesnt *(isnt, mean)* | beauchaine *(bobbeviceicotekcom, away)* |
| 3369591 *(circumference, wk)* | lot *(big, little)* | ive *(seen, never)* |
| mcwilliams *(circumference, wk)* | thats *(unless, youre)* | surrender *(soon, gebcadredslpittedu)* |
| coldblooded *(dictatorship, czech)* | believe *(actually, truth)* | problem *(problems, fix)* |
| militia *(federalist, occupying)* | youre *(unless, theyre)* | windows *(31, dos)* |
| cbc *(ahl, somalia)* | send *(mail, mailing)* | gordon *(soon, gebcadredslpittedu)* |

**Table 3.11: Comparing an MDS sapce to a D2V space for Newsgroups, where a D2V space performed best..**

*Chapter 4*

# Fine-tuning Vector Spaces to Improve Their Directions

"Commonly, these representations are made in a single vector space with similarity being the main structure of interest. However, recent work by Mikolov et al. (2013b) on a word-analogy task suggests that such spaces may have further use- ful internal regularities. They found that seman- tic differences, such as between big and small, and also syntactic differences, as between big and bigger, were encoded consistently across their space. In particular, they solved the word-analogy problems by exploiting the fact that equivalent re- lations tended to correspond to parallel vector- differences. [3]

[3] "Explicitly designing such structure into a neural network model results in rep- resentations that decompose into orthog- onal semantic and syntactic subspaces. We demonstrate that using word-order and morphological structure within En- glish Wikipedia text to enable this decomposition can produce substantial im- provements on semantic-similarity, pos- induction and word-analogy tasks."

This means that despite state-of-the-art results in Natural Language Processing tasks like Language Modelling, Machine Translation, Text Classification, Natural Language Inference, Abstractive Summarization, and Dependency Parsing being dominated by neural networks that learn and improve these kind-of representations, it is not clear what information has been represented.

# 4.1 Experiments

We find that non-linearity is useful.

*Chapter 5*

# Investigating Neural Networks In Terms Of Directions

## 5.1 Appendix

### 5.1.1 Chapter 3 Space Types

| Movies | Genres | | | Keywords | | | Ratings | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | D1 | D2 | D3 | D1 | D2 | D3 | D1 | D2 | D3 |
| Space | 50 PCA | 50 MDS | 100 MDS | 200 PCA | 200 MDS | 200 MDS | 50 PCA | 200 PCA | 50 PCA |
| Single directions | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A |

| Newsgroups | Newsgroups | | | Sentiment | | | Reuters | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | D1 | D2 | D3 | D1 | D2 | D3 | D1 | D2 | D3 |
| Rep | 200 PCA | 200 PCA | 100 PCA | PCA 100 | PCA 50 | PCA 50 | 200 PCA | 200 PCA | 100 PCA |
| Single dir | 200 MDS | 100 D2V | 50 D2V | D2V 100 | PCA 50 | D2V 100 | N/A | N/A | N/A |

| Foursquare | Foursquare | | | OpenCYC | | | Geonames | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | D1 | D2 | D3 | D1 | D2 | D3 | D1 | D2 | D3 |
| Placetypes | MDS 100 | AWV 50 | MDS 200 | AWV 50 | MDS 200 | AWV 50 | MDS 50 | MDS 50 | AWV 200 |
| Single dir | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A | N/A |

**Table 5.1: Space-types, clusters have the same as single directions.**

# GNU Free Documentation License

Version 1.2, November 2002

Copyright © 2000, 2001, 2002 Free Software Foundation, Inc.
59 Temple Place, Suite 330, Boston, MA 02111-1307 USA

Everyone is permitted to copy and distribute verbatim copies of this license document, but changing it is not allowed.

## 0. Preamble

The purpose of this License is to make a manual, textbook, or other functional and useful document *free* in the sense of freedom: to assure everyone the effective freedom to copy and redistribute it, with or without modifying it, either commercially or noncommercially. Secondarily, this License preserves for the author and publisher a way to get credit for their work, while not being considered responsible for modifications made by others.

This License is a kind of "copyleft", which means that derivative works of the document must themselves be free in the same sense. It complements the GNU General Public License, which is a copyleft license designed for free software.

We have designed this License in order to use it for manuals for free software, because free software needs free documentation: a free program should come with manuals providing the same freedoms that the software does. But this License is not limited to software manuals; it can be used for any textual work, regardless of subject matter or whether it is published as a printed book. We recommend this License principally for works whose purpose is instruction or reference.

## 1. Applicability and Definitions

This License applies to any manual or other work, in any medium, that contains a notice placed by the copyright holder saying it can be distributed under the terms of this License. Such a notice grants a world-wide, royalty-free license, unlimited in duration, to use that work under the conditions stated herein. The "Document", below, refers to any such manual or work. Any member of the public is a licensee, and is addressed as "you". You accept the license if you copy, modify or distribute the work in a way requiring permission under copyright law.

A "Modified Version" of the Document means any work containing the Document or a portion of it, either copied verbatim, or with modifications and/or translated into another language.

A "Secondary Section" is a named appendix or a front-matter section of the Document that deals exclusively with the relationship of the publishers or authors of the Document to the Document's overall subject (or to related matters) and contains nothing that could fall directly within that overall subject. (Thus, if the Document is in part a textbook of mathematics, a Secondary Section may not explain any mathematics.) The relationship could be a matter of historical connection with the subject or with related matters, or of legal, commercial, philosophical, ethical or political position regarding them.

The "Invariant Sections" are certain Secondary Sections whose titles are designated, as being those of Invariant Sections, in the notice that says that the Document is released under this License. If a section does not fit the above definition of Secondary then it is not allowed to be designated as Invariant. The Document may contain zero Invariant Sections. If the Document does not identify any Invariant Sections then there are none.

The "Cover Texts" are certain short passages of text that are listed, as Front-Cover Texts or Back-Cover Texts, in the notice that says that the Document is released under this License. A Front-Cover Text may be at most 5 words, and a Back-Cover Text may be at most 25 words.

A "Transparent" copy of the Document means a machine-readable copy, represented in a format whose specification is available to the general public, that is suitable for revising the document straightforwardly with generic text editors or (for images composed of pixels) generic paint programs or (for drawings) some widely available drawing editor, and that is suitable for input to text formatters or for automatic translation to a variety of formats suitable for input to text formatters. A copy made in an otherwise Transparent file format whose markup, or absence of markup, has been arranged to thwart or discourage subsequent modification by readers is not Transparent. An image format is not Transparent if used for any substantial amount of text. A copy that is not "Transparent" is called "Opaque".

Examples of suitable formats for Transparent copies include plain ASCII without markup, Texinfo input format, LaTeX input format, SGML or XML using a publicly available DTD, and standard-conforming simple HTML, PostScript or PDF designed for human modification. Examples of transparent image formats include PNG, XCF and JPG. Opaque formats include proprietary formats that can be read and edited only by proprietary word processors, SGML or XML for which the DTD and/or processing tools are not generally available, and the machine-generated HTML, PostScript or PDF produced by some word processors for output purposes only.

The "Title Page" means, for a printed book, the title page itself, plus such following pages as are needed to hold, legibly, the material this License requires to appear in the title page. For works in formats which do not have any title page as such, "Title Page" means the text near the most prominent appearance of the work's title, preceding the beginning of the body of the text.

A section "Entitled XYZ" means a named subunit of the Document whose title either is precisely XYZ or contains XYZ in parentheses following text that translates XYZ in another language. (Here XYZ stands for a specific section name mentioned below, such as "Acknowledgements", "Dedications", "Endorsements", or "History".) To "Preserve the Title" of such a section when you modify the Document means that it remains a section "Entitled XYZ" according to this definition.

The Document may include Warranty Disclaimers next to the notice which states that this License applies to the Document. These Warranty Disclaimers are considered to be included by reference in this License, but only as regards disclaiming warranties: any other implication that these Warranty Disclaimers may have is void and has no effect on the meaning of this License.

# 2. Verbatim Copying

You may copy and distribute the Document in any medium, either commercially or noncommercially, provided that this License, the copyright notices, and the license notice saying this License applies to the Document are reproduced in all copies, and that you add no other conditions whatsoever to those of this License. You may not use technical measures to obstruct or control the reading or further copying of the copies you make or distribute. However, you may accept compensation in exchange for copies. If you distribute a large enough number of copies you must also follow the conditions in section 3.

You may also lend copies, under the same conditions stated above, and you may publicly display copies.

# 3. Copying in Quantity

If you publish printed copies (or copies in media that commonly have printed covers) of the Document, numbering more than 100, and the Document's license notice requires Cover Texts, you must enclose the copies in covers that carry, clearly and legibly, all these Cover Texts: Front-Cover Texts on the front cover, and Back-Cover Texts on the back cover. Both covers must also clearly and legibly identify you as the publisher of these copies. The front cover must present the full title with all words of the title equally prominent and visible. You may add other material on the covers in addition. Copying with changes limited to the covers, as long as they preserve the title of the Document and satisfy these conditions, can be treated as verbatim copying in other respects.

If the required texts for either cover are too voluminous to fit legibly, you should put the first ones listed (as many as fit reasonably) on the actual cover, and continue the rest onto adjacent pages.

If you publish or distribute Opaque copies of the Document numbering more than 100, you must either include a machine-readable Transparent copy along with each Opaque copy, or state in or with each Opaque copy a computer-network location from which the general network-using public has access to download using public-standard network protocols a complete Transparent copy of the Document, free of added material. If you use the latter option, you must take reasonably prudent steps, when you begin distribution of Opaque copies in quantity, to ensure that this Transparent copy will remain thus accessible at the stated location until at least one year after the last time you distribute an Opaque copy (directly or through your agents or retailers) of that edition to the public.

It is requested, but not required, that you contact the authors of the Document well before redistributing any large number of copies, to give them a chance to provide you with an updated version of the Document.

# 4. Modifications

you may copy and distribute a Modified Version of the Document under the conditions of sections 2 and 3 above, provided that you release the Modified Version under precisely this License, with the Modified Version filling the role of the Document, thus licensing distribution and modification of the Modified Version to whoever possesses a copy of it. In addition, you must do these things in the Modified Version:

**A.** Use in the Title Page (and on the covers, if any) a title distinct from that of the Document, and from those of previous versions (which should, if there were any, be listed in the History section of the Document). You may use the same title as a previous version if the original publisher of that version gives permission.

**B.** List on the Title Page, as authors, one or more persons or entities responsible for authorship of the modifications in the Modified Version, together with at least five of the principal authors of the Document (all of its principal authors, if it has fewer than five), unless they release you from this requirement.

**C.** State on the Title page the name of the publisher of the Modified Version, as the publisher.

**D.** Preserve all the copyright notices of the Document.

**E.** Add an appropriate copyright notice for your modifications adjacent to the other copyright notices.

**F.** Include, immediately after the copyright notices, a license notice giving the public permission to use the Modified Version under the terms of this License, in the form shown in the Addendum below.

**G.** Preserve in that license notice the full lists of Invariant Sections and required Cover Texts given in the Document's license notice.

**H.** Include an unaltered copy of this License.

**I.** Preserve the section Entitled "History", Preserve its Title, and add to it an item stating at least the title, year, new authors, and publisher of the Modified Version as given on the Title Page. If there is no section Entitled "History" in the Document, create one stating the title, year, authors, and publisher of the Document as given on its Title Page, then add an item describing the Modified Version as stated in the previous sentence.

**J.** Preserve the network location, if any, given in the Document for public access to a Transparent copy of the Document, and likewise the network locations given in the Document for previous versions it was based on. These may be placed in the "History" section. You may omit a network location for a work that was published at least four years before the Document itself, or if the original publisher of the version it refers to gives permission.

K. For any section Entitled "Acknowledgements" or "Dedications", Preserve the Title of the section, and preserve in the section all the substance and tone of each of the contributor acknowledgements and/or dedications given therein.

L. Preserve all the Invariant Sections of the Document, unaltered in their text and in their titles. Section numbers or the equivalent are not considered part of the section titles.

M. Delete any section Entitled "Endorsements". Such a section may not be included in the Modified Version.

N. Do not retitle any existing section to be Entitled "Endorsements" or to conflict in title with any Invariant Section.

O. Preserve any Warranty Disclaimers.

If the Modified Version includes new front-matter sections or appendices that qualify as Secondary Sections and contain no material copied from the Document, you may at your option designate some or all of these sections as invariant. To do this, add their titles to the list of Invariant Sections in the Modified Version's license notice. These titles must be distinct from any other section titles.

You may add a section Entitled "Endorsements", provided it contains nothing but endorsements of your Modified Version by various parties — for example, statements of peer review or that the text has been approved by an organization as the authoritative definition of a standard.

You may add a passage of up to five words as a Front-Cover Text, and a passage of up to 25 words as a Back-Cover Text, to the end of the list of Cover Texts in the Modified Version. Only one passage of Front-Cover Text and one of Back-Cover Text may be added by (or through arrangements made by) any one entity. If the Document already includes a cover text for the same cover, previously added by you or by arrangement made by the same entity you are acting on behalf of, you may not add another; but you may replace the old one, on explicit permission from the previous publisher that added the old one.

The author(s) and publisher(s) of the Document do not by this License give permission to use their names for publicity for or to assert or imply endorsement of any Modified Version.

# 5. Combining Documents

You may combine the Document with other documents released under this License, under the terms defined in section 4 above for modified versions, provided that you include in the combination all of the Invariant Sections of all of the original documents, unmodified, and list them all as Invariant Sections of your combined work in its license notice, and that you preserve all their Warranty Disclaimers.

The combined work need only contain one copy of this License, and multiple identical Invariant Sections may be replaced with a single copy. If there are multiple Invariant Sections with the same name but different contents, make the title of each such section unique by adding at the end of it, in parentheses, the name of the original author or publisher of that section if known,

or else a unique number. Make the same adjustment to the section titles in the list of Invariant Sections in the license notice of the combined work.

In the combination, you must combine any sections Entitled "History" in the various original documents, forming one section Entitled "History"; likewise combine any sections Entitled "Acknowledgements", and any sections Entitled "Dedications". You must delete all sections Entitled "Endorsements."

# 6. Collections of Documents

You may make a collection consisting of the Document and other documents released under this License, and replace the individual copies of this License in the various documents with a single copy that is included in the collection, provided that you follow the rules of this License for verbatim copying of each of the documents in all other respects.

You may extract a single document from such a collection, and distribute it individually under this License, provided you insert a copy of this License into the extracted document, and follow this License in all other respects regarding verbatim copying of that document.

# 7. Aggregation with Independent Works

A compilation of the Document or its derivatives with other separate and independent documents or works, in or on a volume of a storage or distribution medium, is called an "aggregate" if the copyright resulting from the compilation is not used to limit the legal rights of the compilation's users beyond what the individual works permit. When the Document is included in an aggregate, this License does not apply to the other works in the aggregate which are not themselves derivative works of the Document.

If the Cover Text requirement of section 3 is applicable to these copies of the Document, then if the Document is less than one half of the entire aggregate, the Document's Cover Texts may be placed on covers that bracket the Document within the aggregate, or the electronic equivalent of covers if the Document is in electronic form. Otherwise they must appear on printed covers that bracket the whole aggregate.

# 8. Translation

Translation is considered a kind of modification, so you may distribute translations of the Document under the terms of section 4. Replacing Invariant Sections with translations requires special permission from their copyright holders, but you may include translations of some or all Invariant Sections in addition to the original versions of these Invariant Sections. You may include a translation of this License, and all the license notices in the Document, and any Warranty Disclaimers, provided that you also include the original English version of this License and the original versions of those notices and disclaimers. In case of a disagreement between

the translation and the original version of this License or a notice or disclaimer, the original version will prevail.

If a section in the Document is Entitled "Acknowledgements", "Dedications", or "History", the requirement (section 4) to Preserve its Title (section 1) will typically require changing the actual title.

# 9. Termination

You may not copy, modify, sublicense, or distribute the Document except as expressly provided for under this License. Any other attempt to copy, modify, sublicense or distribute the Document is void, and will automatically terminate your rights under this License. However, parties who have received copies, or rights, from you under this License will not have their licenses terminated so long as such parties remain in full compliance.

# 10. Future Revisions of this License

The Free Software Foundation may publish new, revised versions of the GNU Free Documentation License from time to time. Such new versions will be similar in spirit to the present version, but may differ in detail to address new problems or concerns. See `http://www.gnu.org/copyleft/`.

Each version of the License is given a distinguishing version number. If the Document specifies that a particular numbered version of this License "or any later version" applies to it, you have the option of following the terms and conditions either of that specified version or of any later version that has been published (not as a draft) by the Free Software Foundation. If the Document does not specify a version number of this License, you may choose any version ever published (not as a draft) by the Free Software Foundation.

# ADDENDUM: How to use this License for your documents

To use this License in a document you have written, include a copy of the License in the document and put the following copyright and license notices just after the title page:

> Copyright © YEAR YOUR NAME. Permission is granted to copy, distribute and/or modify this document under the terms of the GNU Free Documentation License, Version 1.2 or any later version published by the Free Software Foundation; with no Invariant Sections, no Front-Cover Texts, and no Back-Cover Texts. A copy of the license is included in the section entitled "GNU Free Documentation License".

If you have Invariant Sections, Front-Cover Texts and Back-Cover Texts, replace the "with...Texts." line with this:

with the Invariant Sections being `LIST THEIR TITLES`, with the Front-Cover Texts being `LIST`, and with the Back-Cover Texts being `LIST`.

If you have Invariant Sections without Cover Texts, or some other combination of the three, merge those two alternatives to suit the situation.

If your document contains nontrivial examples of program code, we recommend releasing these examples in parallel under your choice of free software license, such as the GNU General Public License, to permit their use in free software.

# Bibliography

[1] Joaqu??n Derrac and Steven Schockaert. Inducing semantic relations from conceptual spaces: A data-driven approach to plausible reasoning. *Artificial Intelligence*, 228:66–94, 2015.

[2] Johan Huysmans, Karel Dejaeger, Christophe Mues, Jan Vanthienen, and Bart Baesens. An empirical evaluation of the comprehensibility of decision table, tree and rule based predictive models. *Decision Support Systems*, 51(1):141–154, 2011.

[3] Jeff Mitchell and Mark Steedman. Orthogonality of Syntax and Semantics within Distributional Spaces. *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, pages 1301–1310, 2015.

[4] Menaka Narayanan, Emily Chen, Jeffrey He, Been Kim, Sam Gershman, and Finale Doshi-Velez. How do Humans Understand Explanations from Machine Learning Systems? An Evaluation of the Human-Interpretability of Explanation. pages 1–21, 2018.

[5] T.L. Saaty and M.S. Ozdemir. Why the magic number seven plus or minus two. *Mathematical and Computer Modelling*, 38(3):233–244, 2003.

[6] Geoffrey Zweig Tomas Mikolovâ , Wen-tau Yih. Linguistic Regularities in Continuous Space Word Representations. *Hlt-Naacl*, (June):746–751, 2013.