

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

Copyright © 2011 Name Lastname.

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”.

A copy of this document in various transparent and opaque machine-readable formats and related software is available at <http://yourwebsite>.

**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

Abstract	iv
Acknowledgements	v
Contents	vi
List of Publications	ix
List of Figures	x
List of Tables	xii
List of Algorithms	xv
List of Acronyms	xvi
0.0.1 Definitions	xvi
1 Introduction	1
1.1 Motivation	1
1.1.1 Directions	3

1.2	Interpretability	4
1.3	Thesis Overview / Contributions	5
2	Background	6
2.1	Text Representations	6
2.1.1	Bag-of-words	6
2.2	Text classification	7
2.2.1	Decision Trees	7
2.2.2	Support Vector Machines	7
2.2.3	Neural Networks	7
2.2.4	Semantic Spaces	7
2.2.5	Document Representations	8
2.3	Interpretable Representations	9
3	Converting Vector Spaces into Interpretable Representations	10
3.1	Introduction	10
3.1.1	Semantic Relations	10
3.1.2	Producing an Interpretable Property Representation	12
3.1.3	Using a Property Representation in a Linear Classifier	13
3.2	Related Work	15
3.2.1	Semantic Relations & Their Applications	15
3.2.2	Interpretable Representations	15
3.3	Method	16

3.3.1	Term Rankings	17
3.3.2	Filtering Words	19
3.3.3	Qualitative Results	22
3.3.4	Quantitative Results	22
3.3.5	Interpretability Results	26
4	Fine-tuning Vector Spaces to Improve Their Directions	27
4.1	Experiments	27
5	Investigating Neural Networks In Terms Of Directions	28
	GNU Free Documentation License	29
	Bibliography	38

List of Publications

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

-
-

List of Figures

1.1	Bag-of-words	2
1.2	Example properties	2
3.1	Movies and selected associated dimensions, and their use in a linear classifier.	11
3.2	This figure shows a 2d toy space where entities are shapes and directions are properties. We demonstrate on the right the method to induce a ranking from the directions, in particular by using the dot-product of the entity point on the directions vector. In the same way for a more complex space, we can understand each entity point to be ranked on thousands of property directions, and the space to be much higher dimensionality.	13
3.3	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.	14
3.4	Original And Converted.	17
3.5	Original And Converted.	18

3.6	A conceptual space of movies, where regions correspond to properties and entities are points.	23
-----	--	----

List of Tables

- 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. 24

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 w for its document $D_w f$.

wf

Bag-Of-Words a matrix BOW of documents BOW_D where each document is composed of unordered frequencies of words $D = [wf_1, \dots, wf_n]$. 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

<u>Entity: X</u>		<u>Entity: Y</u>		<u>Entity: Z</u>	
<u>Word</u>	<u>Frequency</u>	<u>Word</u>	<u>Frequency</u>	<u>Word</u>	<u>Frequency</u>
Dog	51	Dog	51	Dog	51
Cat	40	Cat	40	Cat	40
Man	11	Man	11	Man	11
Cheese	0	Cheese	0	Cheese	0
Dog	51	Dog	51	Dog	51
Cat	40	Cat	40	Cat	40
Man	11	Man	11	Man	11
Cheese	0	Cheese	0	Cheese	0

Figure 1.1: Bag-of-words

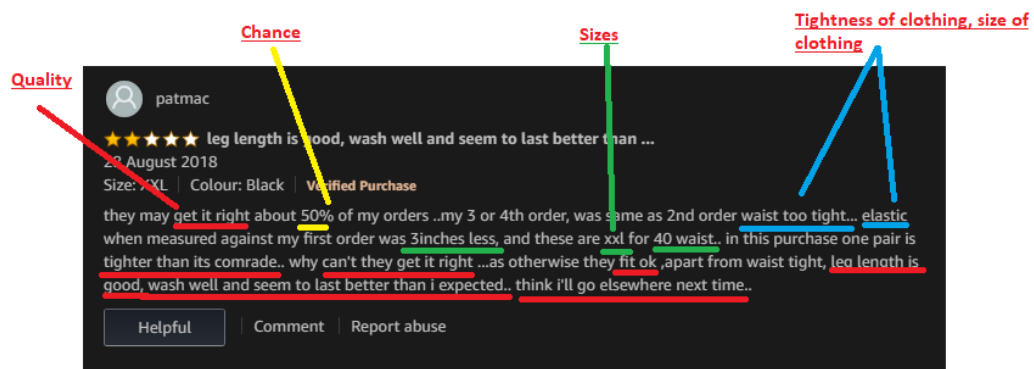


Figure 1.2: Example properties

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 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 shiner”, 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' [10] also has many positive effects for its users, like lower response times [9, 7], better question answering and confidence for logical problem questions [7] and higher satisfaction [9].
- Transparency:
- 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 [6]: 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.

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 lowercase 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

2.2.2 Support Vector Machines

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

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.

Chapter 3

Converting Vector Spaces into Interpretable Representations

3.1 Introduction

Semantic spaces encode similarity information spatially, where the spatial position of entities in the space is affected by its semantic relationship to the other entities. Semantic Spaces are able to accurately represent complex domain meaning, achieving state-of-the-art results in X, Y, Z. However, the dimensions of these spaces do not correspond to the kind-of features we typically understand to describe entities, for example describing a movie by its genre. Additionally, the dimensions of the space and the relationships between entities are not labelled with words. This makes it difficult to apply in X, Y, Z. Although the dimensions do not correspond to features, the space does accurately represent these features through spatial relationships. In this section, we show how to go from a bag-of-words and standard Semantic Space to a representation where each dimension corresponds to a ranking of the entity on a labelled feature of the domain, e.g. a movie would be ranked on dimensions like 'Horror, Gore, Scary', and 'Romance, Love, Relationships'. We then show that these features perform better than Semantic Spaces and a baseline Topic Model when applied to a simple linear classifiers (Decision Tree) and are more semantically cohesive and interpretable than a Topic Model. We show an example representation and classifier that were obtained

and rep.png

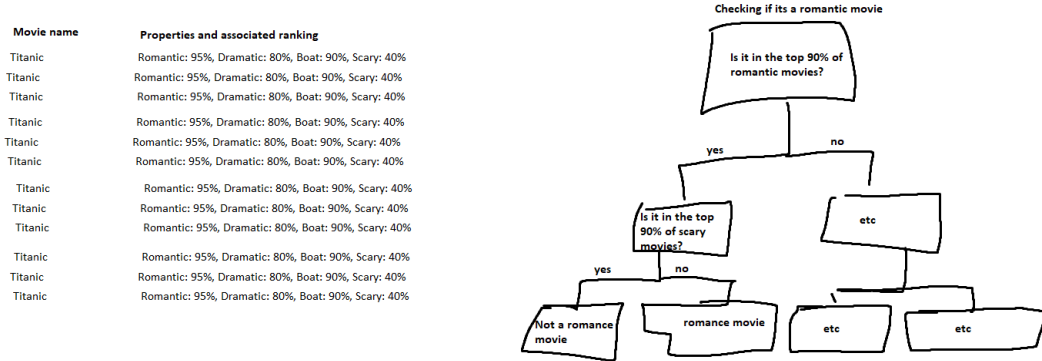


Figure 3.1: Movies and selected associated dimensions, and their use in a linear classifier..

from a Semantic Space constructed from IMDB movie reviews in figure 3.1.

3.1.1 Semantic Relations

The success of these semantic spaces similarity-based structure has lead many to investigate how to formalize the relationships they encode. One example is in that of linear analogies in word-vectors (see Section ??, where it was found that the vector $XXXX[King queen blah blah[XXXXXXX$, formally justified in [?]. These relationships have been expanded on, for example [11] found that "equivalent relations tended to correspond to parallel vector differences" [8], while [8] discovered that by decomposing representations into orthogonal semantic and syntactic subspaces they were able to produce substantial improvements on various tasks. Additionally, [?] found that 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.

In this paper, we use directions in the space that correspond to salient features from the considered domain. A "direction" refers in this case to the orthogonal direction to a hyper plane that separates a term in a vector space. As the hyper plane separates

entities, this means that the entities furthest along the hyper plane, at the end classified positively, are the entities we are most sure have the term we found the hyper plane for. To see an example of this, see ?? With this understanding, it becomes possible to induce a ranking of entities on the properties by finding the dot product of the entity points on the direction vector. These kind-of directions have been used in many different ways for different domains, For instance, [?] found that features of countries, such as their GDP, fertility rate or even level of CO₂ 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.

3.1.2 Producing an Interpretable Property Representation

By finding the dot product between entity points in the space and direction vectors, it is possible to induce a ranking of entities on those directions. In this chapter, we more deeply investigate the potential of direction vectors to rank entities on properties to form an interpretable representation. In this thesis, we refer to these direction vectors as directions to convey the ordinal meaning, and directions as feature-directions if they are sufficiently salient in the space, e.g. In a domain of IMDB movie reviews where movies are entities, a direction on the word "The" would not be a feature-direction, but a direction on the word "Horror" would be.

We demonstrate the effect of different filtering methods to find properties, the ability of different clustering methods to label properties, as well as the number and types of directions, for use in a Decision Tree. In Figure 3.3, we demonstrate how depth could affect a Decision Tree that uses salient feature-directions. Decision Tree's have the additional benefit of giving us a way to identify the degree to which our feature-directions are salient - if we are able to construct a simple but high-scoring Decision Tree limited to only a single node for if a movie has the genre of 'Comedy' using only our ranking

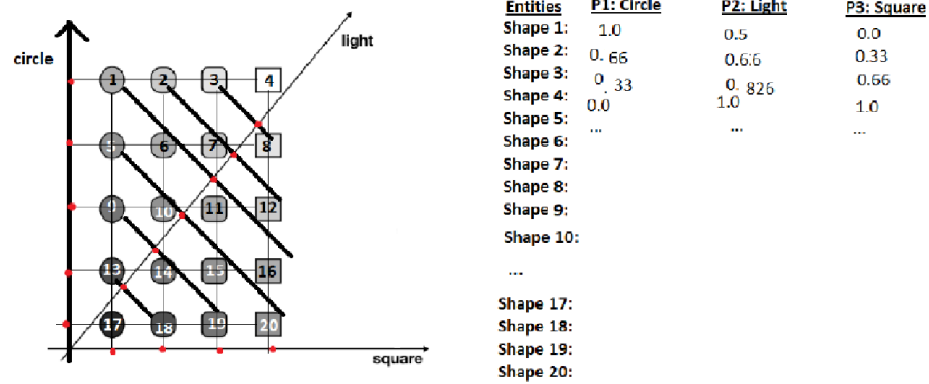


Figure 3.2: This figure shows a 2d toy space where entities are shapes and directions are properties. We demonstrate on the right the method to induce a ranking from the directions, in particular by using the dot-product of the entity point on the directions vector. In the same way for a more complex space, we can understand each entity point to be ranked on thousands of property directions, and the space to be much higher dimensionality..

of entities on the feature-direction $p = \text{"Funny"}, \text{"Hilarious"}, \text{"Laughing"}$ then we know that this feature-direction is salient.

This chapter continues as follows: We begin by describing the work related to this , giving valuable context for the utility and potential of our approach. This is followed by an explanation of the method, including the variations we have adopted for our experimental work. We follow this with our qualitative experimentation, explaining how these variations affect the results, as well as the interpretability of the method, and we end with a quantitative analysis on how well we can represent domain knowledge using decision trees constrained to a limited depth.

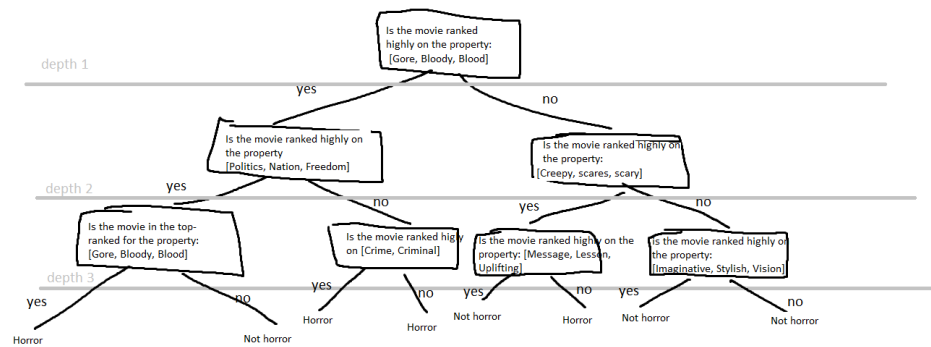


Figure 3.3: This figure shows an example tree from one of our classifiers. Here, we can see that the model increases in complexity as it increases in depth. In this case, we end-up getting better F-score with just a depth-one tree, as the tree begins to overfit at depth three. .

3.2 Related Work

3.2.1 Semantic Relations & Their Applications

Linear Classifiers Decision trees, linear SVM's, logistic regression, decision tables, IF Then rules.

What are the available options for interpretable linear classification?

How have each of these methods been measured or validated in the literature in regards to interpretability? How about application to real world situations?

Non linear classifiers What non linear classifiers networks are interpretable? How have they done it? How have they measured it? How does it compare to a linear method?

Neural networks Approximating w/linear model, Interpretable nodes/weights

Other Stuff

3.2.2 Interpretable Representations

There are two ways in which topic models can be used for document classification. First, a supervised topic model can be used, in which the underlying graphical model is explicitly extended with a variable that represents the class label [4]. Second, the parameters of the multinomial distribution corresponding to a given document can be used as a feature vector for a standard classifier, such as a Support Vector Machine (SVM) or Decision Tree. LDA has been extended by many approaches, e.g. aiming to avoid the need to manually specify the number of topics [?], modelling correlations between topics [3], or by incorporating meta-data such as authors [?] or time stamps [?].

Broadly speaking, in the context of document classification, the main advantage of topic models is that their topics tend to be easily interpretable, while vector space models tend to be more flexible in the kind of meta-data that can be exploited. The approach we propose in this paper aims to combine the best of both worlds, by providing a way to derive interpretable representations from vector space models.

3.3 Method

This section details the methodology to go from a Bag-Of-Words (BOW) 2.1.1 and Semantic Space 2.2.4, to interpretable vectors that rank documents on features of the domain, e.g. A movie would be ranked on how *Scary*, *Horror*, *Bloody* it is for one dimension of the feature-vector, and how *Romantic*, *Love*, *Cute* it is in another, ideally with as many dimensions as there are distinct salient features of the domain. We show examples of this final representation in ?? . For the Bag-Of-Words, we begin with an unstructured corpus of text documents from a domain, e.g. movie reviews, where each document is a collection of reviews for a movie. From these reviews, we preprocess the text such that it is converted to lower-case, and non-alphanumeric characters are

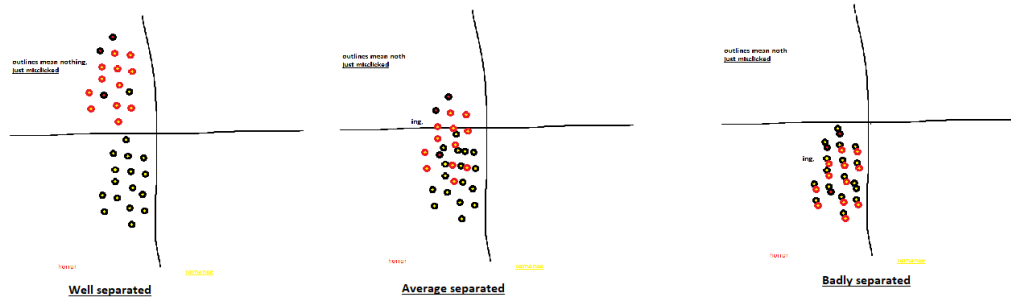


Figure 3.4: Original And Converted.

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 3.4. From this preprocessed corpus, we obtain a Bag-Of-Words where we count the frequency of each term BOW_{wf} , see 2.1.1. 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 ??.

The method to obtain interpretable feature-vectors is an extension of the work by [?]. This previous work showed how to, filter out words, cluster words to get features, and obtain rankings of documents on those features. In this section, we further analyse and extend this work, in particular by testing a variety of additional filtering methods and clustering methods, and demonstrating how these feature-vectors can produce simple linear interpretable classifiers.

3.3.1 Term Rankings

Structure of a Semantic Space Salient features of the domain are encoded in the structure of a semantic space, see Section 2.2.4 for more detail. We can expect that for these salient features, they will be more linearly separable than words, and be spatially organized in a way that reflects the similarity between their associated PPMI scores

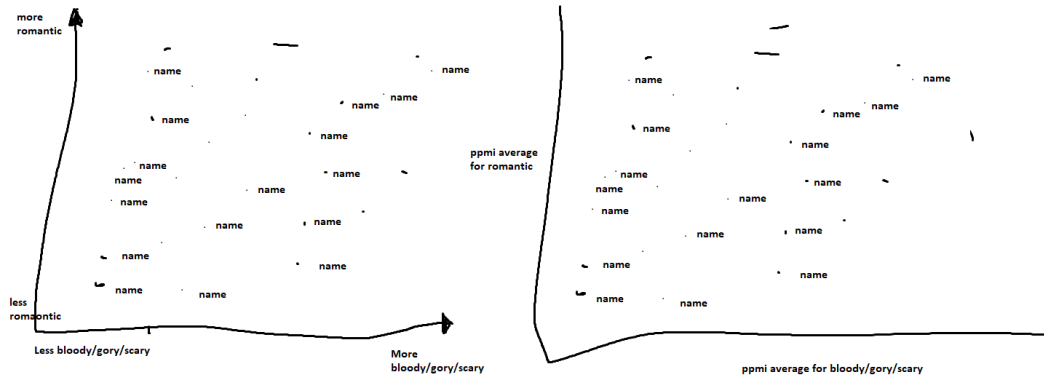


Figure 3.5: Original And Converted.

that the space was constructed from. In particular, we expect that documents will be arranged in a direction, where generally the higher the PPMI score for a group of words that correspond to a feature (e.g. *Horror*, *Scary*, *Gore*) the further away they will be from those that have low PPMI scores for those words. We give examples of this in Figure 3.5, by projecting documents into a 2D space of salient features we are able to show that these documents are structured according to directions for these features. Salient features will typically be a more abstract representation which will be natural in the domain, e.g. in a domain of IMDB movie reviews, genres. However, in this section we show how to extract rankings of documents on words, with the understanding that all words may not be features of the domain. In the next sections, we aim to use these words to extract salient features by filtering and clustering.

Obtaining directions for each word For each word w , a Support Vector Machine (See Section ??) 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 $w_f \geq 1$ and negative otherwise. Although the separation of documents is binary, given the structure of the semantic space we can expect for salient features that the documents close to the hyper-plane on the positive side will have lower PPMI scores for the term

than those furthest from the hyper-plane on the positive side, as they are closest to the documents that are classified negatively. 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. In ??, we show an example of directions in a toy domain.

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 an 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_{dw} of the document d for the word w , and in particular, we take $r_{d1} < r_{d2}$ to mean that d_2 has the property labelled with the word w to a greater extent than d_1 .

3.3.2 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_{dw} . 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 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, X proposed to use Cohen’s Kappa score instead. In our experiments, however, we found that accuracy sometimes yields better results, so 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 results in this more abstract feature. As 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.

The final benefit to clustering the words is that linear classifiers are generally suited better to 'disentangled' representations [1]. In this case, we refer to disentanglement in the sense of obtaining a feature vector where each dimension is distinct, rather than the semantic space being naturally clustered. Additionally, if our representation is dense and disentangled into the natural features of the domain, it is unlikely to overfit and will be able to generalize more easily. When investigating the use of directions without clustering in Section ?? we found that the sparsity of the directions when using only words tends to overfit in simple linear classifiers.

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, \dots, 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}, \dots, v_{C_k} are the feature directions that are identified by our method.

3.3.3 Qualitative Results

Dimension of the Space

Space-type

Scoring method

Clustering method

The effect of dimensions

The effect of space-type

The effect of scoring method

3.3.4 Quantitative Results

Datasets

Vocabulary size/entity size/origin of data/classes

Semantic Spaces

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

As the newsgroups contained empty documents after removing all words that do not occur in at least 2 documents, we have removed these empty documents, leaving us with 18302 overall documents. Following this, instead of using the train split as determined by previous literature, we did a simple 2/3 train/test split the same as our IMDB dataset.

This section focuses on using linear classifiers to determine how well our method represents domain knowledge compared to standard baselines. We can understand that an accurate representation of domain knowledge will be one that ensures semantically distinct entities are separated, and semantically similar entities are close together. Put another way, if the space is representing domain knowledge well we can expect that the space should be linearly separable for key semantics of the domain. 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 3.6. For a Bag-Of-Words, we can expect similar entities to have similarly scoring terms ??.

When selecting the parameters to use for the doc2vec space when obtaining directions, we choose the one that scored the highest for its class on a Linear SVM, rather than

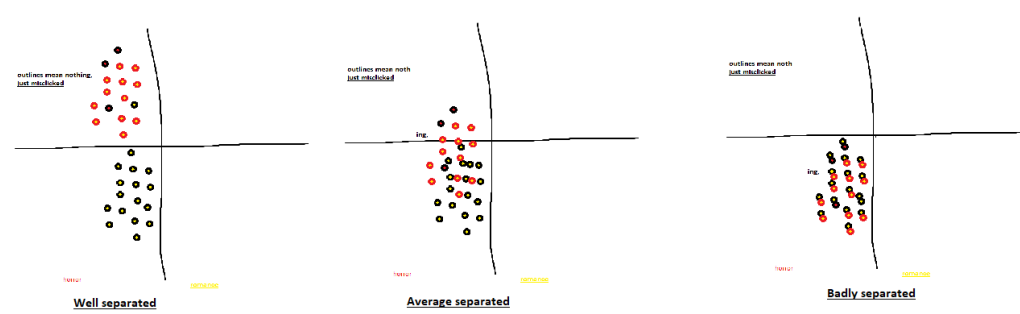


Figure 3.6: A conceptual space of movies, where regions correspond to properties and entities are points..

	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..

tuning the entire process around the doc2vecs vectors. We are not able to obtain an MDS space for sentiment or doc2vec spaces for placetypes/movies.

We obtain results with just these spaces as input, and additionally results for a bag-of-words PPMI representation. These results act as our baselines for quantitative results, in addition to a Topic Model. We find results using a Linear SVM, and a Decision Tree with an unlimited depth, a depth of 3, a depth of 2, and a depth of 1. Each SVM is tuned using a grid search for the optimal C value, and whether or not to balance the classes. For all trees we attempt to find the best value between [None, 'auto', 'log2'],

and additionally try differnet criterion, either the gini score or the information entropy score. In the same way as the SVM's, we include whether or not to balance the classes in the grid search.

Word Directions

The binary BOW representation for each word that has not been removed by the frequency cut-off is used as a target for a linear SVM, with a Semantic Space as features. We use the scikit-learn libraries LinearSVC implementation with a default C value (1.0). We balance the classes, as many of the binary BOW representations are sparse, and use the primal formulation.

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, 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. 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.

What is the importance of the space size? What is the importance of the space type? How do directions perform compared to spaces? Why? What kind of directions do we find for each score type? Why does the score type matter? What score type works best?

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. These single directions typically overfit.

Clustered Directions

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 infrequent 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-learn¹, introduced by [?] and the kmeans++

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

to initialize [?]

3.3.5 Interpretability Results

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 useful internal regularities. They found that semantic 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 relations tended to correspond to parallel vector-differences. [8]

[8] "Explicitly designing such structure into a neural network model results in representations that decompose into orthogonal semantic and syntactic subspaces. We demonstrate that using word-order and morphological structure within English Wikipedia text to enable this decomposition can produce substantial improvements on semantic-similarity, pos-induction and word-analogy tasks."

4.1 Experiments

We find that non-linearity is useful.

Chapter 5

Investigating Neural Networks In Terms Of 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 non-commercially. Secondly, 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] Yoshua Bengio, Aaron Courville, and Pascal Vincent. Representation Learning: A Review and New Perspectives. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 35(8):1798–1828, 2012.
- [2] David M Blei, Blei@cs Berkeley Edu, Andrew Y Ng, Ang@cs Stanford Edu, Michael I Jordan, and Jordan@cs Berkeley Edu. Latent Dirichlet Allocation. *Journal of Machine Learning Research*, 3:993–1022, 2003.
- [3] David M. Blei and John D. Lafferty. Correlated Topic Models. *Advances in Neural Information Processing Systems 18*, pages 147–154, 2006.
- [4] David M. Blei and Jon D. McAuliffe. Supervised Topic Models. pages 1–8, 2010.
- [5] Jonathan Chang, Sean Gerrish, Chong Wang, and David M Blei. Reading Tea Leaves: How Humans Interpret Topic Models. *Advances in Neural Information Processing Systems 22*, pages 288—296, 2009.
- [6] Joaquin Derrac and Steven Schockaert. Inducing semantic relations from conceptual spaces: A data-driven approach to plausible reasoning. *Artificial Intelligence*, 228:66–94, 2015.
- [7] 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.
- [8] 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.

-
- [9] 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.
- [10] 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.
- [11] Geoffrey Zweig Tomas Mikolov , Wen-tau Yih. Linguistic Regularities in Continuous Space Word Representations. *Hlt-Naacl*, (June):746–751, 2013.