

# Application

---

1. Application Number: 3
2. Post Applied for: Associate Professor grade2
3. Name: Shreya B Ballijepalli
4. Category: General
5. Disability: NO
6. Date of Birth: 2019-05-10
7. Nationality: India
8. Gender: Female
9. Marital status: Married
10. Address: 1-19-80/103B, Vijayapuri Colony
11. Email ID: shreyabalijepalli@gmail.com

## 12. Education

Exam Passed	Board/university	Year of Passing	Specialization	CGPA/Percentage
Bachelors	de	2019-05-09	ss	4
Masters	rf	2019-05-22	3	1

## 13. PHD:

University	Year of Graduation	Date of thesis submission	Date of Defence	Specialization	CGPA
hhe	2019-05-11	2019-05-21	2019-05-22	swd	1

## 14. GATE Year: 1999

## 15. GATE Score: 1

16. Research Specialization: (2019-05-21,2019-05-22)

17. Research Interests: a

18. Post Doc Specialization: a

19. Present Position with Salary Details:

Position	Pay Band	Grade Pay	Consolidated Salary
----------	----------	-----------	---------------------

20. Research/Teaching/Industrial Experience(if any):

Name of the Organization	Start Date	End Date	Full Time (Yes/No)	Designation	Type of Work
--------------------------	------------	----------	--------------------	-------------	--------------

21. Projects

Type of Project	project Title	Project Amount	Project Details
-----------------	---------------	----------------	-----------------

22. Referees

Name	Email	Designation	Address
Shreya Ballijepalli	cs15btech11009@ii th.ac.in	1	1-19-80/103B
"Shreya Ballijepalli"	cs15btech11009@ii th.ac.in	wer	1
erty	cs15btech11009@ii th.ac.in	1	1



Please use `late` instead of `keep_prev`, `late` should be set to `late - 1 - keep_prev`.

Layer (type)	Output Shape	Param #
dense_1 (Dense)	(None, 100)	29900
batch_normalization_1 (Batch Normalization)	(None, 100)	400
activation_1 (Activation)	(None, 100)	0
dropout_1 (Dropout)	(None, 100)	0
dense_2 (Dense)	(None, 80)	8080
batch_normalization_2 (Batch Normalization)	(None, 80)	320
activation_2 (Activation)	(None, 80)	0
dropout_2 (Dropout)	(None, 80)	0
batch_normalization_3 (Batch Normalization)	(None, 80)	320
activation_3 (Activation)	(None, 80)	0
dropout_3 (Dropout)	(None, 80)	0
dense_3 (Dense)	(None, 9)	729
batch_normalization_4 (Batch Normalization)	(None, 9)	36
activation_4 (Activation)	(None, 9)	0
Total params: 39,785		
Trainable params: 39,247		
Non-trainable params: 538		

```
In [6]: print(y_train)
print(len(x_train), y_train.shape)
```

## Case Study - Walmart

B.Shreya  
cs15btech11009

Walmart is an American multinational retail corporation that was started in 1962 by Sam Walton. It operates a chain of hypermarkets, department stores, grocery units and it is the largest private employer in the US as well as the world's largest retailer. Its corporate mission is "**Save people money so they can live better**". It is known as the "global legislator" because it's an important emerging private actor in the transformation of lawmaking in the CSR field.

### **Effect of Walmart's Expansion:**

Walmart was first started in 1962 with its target towards rural towns with a population of lesser than ten thousand people. It then extended the company to large cities and opened international stores across the world.

#### Affects on businesses of small local merchants:

Whenever Walmart starts a new store in a town, all the small merchants and business vendors feel uncomfortable and fear that they can no longer compete with low prices offered by Walmart. There have also been cases where some merchants have pulled out their businesses when Walmart entered their town. This is also known as "**The Walmart Effect**".

#### Creation of Urban Sprawl, Traffic Congestion:

Walmart's mega stores are built on vast areas of land. By doing this, Walmart is depressing the economic health of communities and other downtown stores. These stores are also built in areas which are not accessible without driving resulting in a lot of traffic.

#### Environmental Pollution:

The increased traffic has led to more air pollution, water contamination and call for more roads. The landscape is also affected because of these stores as they have huge areas with many unused parking lots as well.





### **Walmart's initiative towards Corporate Social Responsibility:**

From 2007, Walmart publishes its annual report on its website which is known as the 'Global Responsibility Report'. This report talks about Walmart's constant and progressive efforts towards social responsibility issues. It has made investments in education, health, commitments to fight hunger, support to local farmers.

### **Walmart's Conflicts:**



### **Gender Discrimination:**

There have been several charges of gender and racial discrimination on Walmart. Walmart Stores vs Duke et Al was one such case filed on Walmart in 2001 and it is the largest class action lawsuit in US history. The plaintiff class included 1.6 million women who were lead by Betty Dukes. Dukes was a 54-year-old Walmart worker in California who claimed that despite six years of work and positive performance reviews, she was denied training she needed to advance to a higher position.

Dukes and others claimed that women were discriminated against pay and promotions to top management positions violating the Civil Rights Act of 1964. Walmart appealed to the Ninth Circuit in 2005 that the seven lead plaintiffs were not typical or common of the class. Walmart then turned to the Supreme Court in 2010 after the Ninth Circuit court upheld class certification. In 2011, the Supreme Court reversed class certification saying that the millions of plaintiffs and their claims didn't have enough in common. This case, however, didn't end here as the plaintiffs filed an amended lawsuit in October 2011, limiting the class to female employees in California.

After filling of the lawsuit, Walmart incorporated an Advisory Board on Gender Equality and Diversity, which is aimed at providing equal opportunities for all in top management positions. It has also included a Gender Equality and Diversity gender Policy in its 'Global Responsibility Annual Report'.

Below is a picture from Walmart's 2015 Diversity an Inclusion Report.

### Diversity Goals Program

Our Diversity Goals Program is the most significant means by which we have accelerated opportunity for our women and people of color associates in the U.S.

The program encompasses:

- Field management placement goals of women and people of color associates
- Good Faith Efforts to drive ownership of diversity and inclusion
- Five-year aspirational goals to stretch our management placement goals for store and club manager positions
- Active coaching reviews centered on discrimination and harassment
- Customized diversity and inclusion plans for senior leaders



### WOMEN REPRESENTATION



As of January 31, 2015

### PEOPLE OF COLOR REPRESENTATION



### Child Labour:

In 2005, a Radio Canada programme Zone Libre reported that Walmart was using child labor at two factories in Bangladesh. Walmart employed children between the ages of 10-15 years for less than \$50 a month for manufacturing products and exporting them to Canada.

After this incident, Walmart ceased businesses with the two factories immediately. In a 2005 Ethical Sourcing Report of Walmart, it stated that Walmart ceased to do business with 141 companies because of underage labor violations. The stakeholders affected in this were thousands of poor workers who lost their jobs as a result of this.



Walmart's 2005 and 2012 COC 'Standard for Suppliers' explicitly establish it would not tolerate the use of child labor and it sets 14 as the minimum age for any worker.

### **Walmart's Bribery Scandal:**



Walmart de Mexico, one of the most successful businesses of Walmart was caught in a massive bribery scandal in April 2012. The bribes which totaled more than \$24 million were given to the Mexican government to win permission to open stores at a much rapid phase. (which wouldn't have been possible according to the Mexican laws). Walmart's senior management long knew about the scandal and tried to cover it up. When this case came into light, it was suggested that Walmart undergo a harsh investigation. However, Walmart opted for an in-house investigation and gave the primary responsibility of the investigation to Walmart de

Mexico itself, again another attempt to conceal the fraud. This was not surprisingly “quickly discontinued.”

Walmart used bribery as a mean to monopolize, neglecting the rules that are set to protect a town and its inhabitants from unsafe commercial development. This action also affected the local businesses to a great extent as consumers were driven towards the “low prices” offered by Walmart.

After this scandal became public, Walmart suffered investor lawsuits, numerous investigations from the Department of Justice and Securities and also brand damage. The stakeholders affected in this scandal were Walmart’s investors, local businesses.

This scandal is still under investigation and it is predicted that criminal charges for some of the Walmart executives are certain.

### **Conclusion:**

Walmart is becoming internationally strong and big day by day. Its low prices have really grabbed customers and have resulted in the shut down of a lot of local businesses. There are several organizations like “Wakeup Walmart”, “Walmart March” which are fighting against the company and the company has also been part of several allegations like poor working conditions, low wages, undertrained workers, etc but it is still going strong day by day.

### **Resources:**

<https://www.scribd.com/document/373615247/walmart-corporate-social-responsibility-case-study>

<https://www.ukessays.com/essays/management/understanding-of-the-case-study-walmart-management-essay.php>

[https://en.wikipedia.org/wiki/Criticism\\_of\\_Walmart](https://en.wikipedia.org/wiki/Criticism_of_Walmart)

<https://www.scribd.com/document/342567422/Case-Study-Over-Csr-Conflicts>

<https://www.businessinsider.com/walmart-bribery-scandal-2012-4?IR=T>

<https://cdn.corporate.walmart.com/01/8b/4e0af18a45f3a043fc85196c2cbe/2015-diversity-and-inclusion-report.pdf>

[https://www.academia.edu/10316637/CASE\\_STUDY\\_MEXICO\\_WALMART\\_SCANDAL](https://www.academia.edu/10316637/CASE_STUDY_MEXICO_WALMART_SCANDAL)

---

# Hierarchical Dirichlet Processes

---

Yee Whye Teh<sup>(1)</sup>, Michael I. Jordan<sup>(1,2)</sup>, Matthew J. Beal<sup>(3)</sup> and David M. Blei<sup>(1)</sup>

<sup>(1)</sup>Computer Science Div., <sup>(2)</sup>Dept. of Statistics

University of California at Berkeley

Berkeley CA 94720, USA

{ywtteh, jordan, blei}@cs.berkeley.edu

<sup>(3)</sup>Dept. of Computer Science

University of Toronto

Toronto M5S 3G4, Canada

beal@cs.toronto.edu

## Abstract

We propose the hierarchical Dirichlet process (HDP), a nonparametric Bayesian model for clustering problems involving multiple groups of data. Each group of data is modeled with a mixture, with the number of components being open-ended and inferred automatically by the model. Further, components can be shared across groups, allowing dependencies across groups to be modeled effectively as well as conferring generalization to new groups. Such grouped clustering problems occur often in practice, e.g. in the problem of topic discovery in document corpora. We report experimental results on three text corpora showing the effective and superior performance of the HDP over previous models.

## 1 Introduction

One of the most significant conceptual and practical tools in the Bayesian paradigm is the notion of a *hierarchical model*. Building on the notion that a parameter is a random variable, hierarchical models have applications to a variety of forms of grouped or relational data and to general problems involving “multi-task learning” or “learning to learn.” A simple and classical example is the Gaussian means problem, in which a grand mean  $\mu_0$  is drawn from some distribution, a set of  $K$  means are then drawn independently from a Gaussian with mean  $\mu_0$ , and data are subsequently drawn independently from  $K$  Gaussian distributions with these means. The posterior distribution based on these data couples the means, such that posterior estimates of the means are shrunk towards each other. The estimates “share statistical strength,” a notion that can be made precise within both the Bayesian and the frequentist paradigms.

Here we consider the application of hierarchical Bayesian ideas to a problem in “multi-task learning” in which the “tasks” are clustering problems, and our goal is to share clusters among multiple, related clustering problems. We are motivated by the task of discovering topics in document corpora [1]. A topic (i.e., a cluster) is a distribution across words while documents are viewed as distributions across topics. We want to discover topics that are common across multiple documents in the same corpus, as well as across multiple corpora.

Our work is based on a tool from nonparametric Bayesian analysis known as the *Dirichlet process* (DP) mixture model [2, 3]. Skirting technical definitions for now, “nonparametric” can be understood simply as implying that the number of clusters is open-ended. Indeed, at each step of generating data points, a DP mixture model can either assign the data point

to a previously-generated cluster or can start a new cluster. The number of clusters is a random variable whose mean grows at rate logarithmic in the number of data points.

Extending the DP mixture model framework to the setting of multiple related clustering problems, we will be able to make the (realistic) assumption that we do not know the number of clusters a priori in any of the problems, nor do we know how clusters should be shared among the problems.

When generating a new cluster, a DP mixture model selects the parameters for the cluster (e.g., in the case of Gaussian mixtures, the mean and covariance matrix) from a distribution  $G_0$ —the *base distribution*. So as to allow any possible parameter value, the distribution  $G_0$  is often assumed to be a smooth distribution (i.e., non-atomic). Unfortunately, if we now wish to extend DP mixtures to groups of clustering problems, the assumption that  $G_0$  is smooth conflicts with the goal of sharing clusters among groups. That is, even if each group shares the same underlying base distribution  $G_0$ , the smoothness of  $G_0$  implies that they will generate distinct cluster parameters (with probability one).

We will show that this problem can be resolved by taking a hierarchical Bayesian approach. We present a notion of a *hierarchical Dirichlet process* (HDP) in which the base distribution  $G_0$  for a set of Dirichlet processes is itself a draw from a Dirichlet process. This turns out to provide an elegant and simple solution to the problem of sharing clusters among multiple clustering problems.

The paper is organized as follows. In Section 2, we provide the basic technical definition of DPs and discuss related representations involving stick-breaking processes and Chinese restaurant processes. Section 3 then introduces the HDP, motivated by the requirement of a more powerful formalism for the grouped data setting. As for the DP, we present analogous stick-breaking and Chinese restaurant representations for the HDP. We present empirical results on a number of text corpora in Section 5, demonstrating various aspects of the HDP including its nonparametric nature, hierarchical nature, and the ease with which the framework can be applied to other realms such as hidden Markov models.

## 2 Dirichlet Processes

In this section we give a brief overview of Dirichlet processes (DPs) and DP mixture models, with an eye towards generalization to HDPs. We begin with the definition of DPs [4]. Let  $(\Theta, \mathcal{B})$  be a measurable space, with  $G_0$  a probability measure on the space, and let  $\alpha_0$  be a positive real number. A *Dirichlet process* is the distribution of a random probability measure  $G$  over  $(\Theta, \mathcal{B})$  such that, for any finite partition  $(A_1, \dots, A_r)$  of  $\Theta$ , the random vector  $(G(A_1), \dots, G(A_r))$  is distributed as a finite-dimensional Dirichlet distribution:

$$(G(A_1), \dots, G(A_r)) \sim \text{Dir}(\alpha_0 G_0(A_1), \dots, \alpha_0 G_0(A_r)) . \quad (1)$$

We write  $G \sim \text{DP}(\alpha_0, G_0)$  if  $G$  is a random probability measure distributed according to a DP. We call  $G_0$  the base measure of  $G$ , and  $\alpha_0$  the concentration parameter.

The DP can be used in the mixture model setting in the following way. Consider a set of data,  $\mathbf{x} = (x_1, \dots, x_n)$ , assumed exchangeable. Given a draw  $G \sim \text{DP}(\alpha_0, G_0)$ , independently draw  $n$  *latent factors* from  $G$ :  $\phi_i \sim G$ . Then, for each  $i = 1, \dots, n$ , draw  $x_i \sim F(\phi_i)$ , for a distribution  $F$ . This setup is referred to as a *DP mixture model*.

If the factors  $\phi_i$  were all distinct, then this setup would yield an (uninteresting) mixture model with  $n$  components. In fact, the DP exhibits an important *clustering property*, such that the draws  $\phi_i$  are generally *not* distinct. Rather, the number of distinct values grows as  $O(\log n)$ , and it is this that defines the random number of mixture components.

There are several perspectives on the DP that help to understand this clustering property. In this paper we will refer to two: the *Chinese restaurant process* (CRP), and the *stick-*



*breaking process.* The CRP is a distribution on partitions that directly captures the clustering of draws from a DP via a metaphor in which customers share tables in a Chinese restaurant [5]. As we will see in Section 4, the CRP refers to properties of the joint distribution of the factors  $\{\phi_i\}$ . The stick-breaking process, on the other hand, refers to properties of  $G$ , and directly reveals its discrete nature [6]. For  $k = 1, 2, \dots$ , let:

$$\theta_k \sim G_0 \quad \beta'_k \sim \text{Beta}(1, \alpha_0) \quad \beta_k = \beta'_k \prod_{l=1}^{k-1} (1 - \beta'_l). \quad (2)$$

Then with probability one the random measure defined by  $G = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}$  is a sample from  $\text{DP}(\alpha_0, G_0)$ . The construction for  $\beta_1, \beta_2, \dots$  in (2) can be understood as taking a stick of unit length, and repeatedly breaking off segments of length  $\beta_k$ . The stick-breaking construction shows that DP mixture models can be viewed as mixture models with a countably infinite number of components. To see this, identify each  $\theta_k$  as the parameter of the  $k^{\text{th}}$  mixture component, with mixing proportion given by  $\beta_k$ .

The DP and the DP mixture are mainstays of nonparametric Bayesian statistics (see, e.g., [3]). They have also begun to be seen in applications in machine learning (e.g., [7, 8, 9]).

### 3 Hierarchical Dirichlet Processes

We will introduce the *hierarchical Dirichlet process* (HDP) in this section. First we describe the general setting in which the HDP is most useful—that of *grouped data*.

We assume that we have  $J$  groups of data, each consisting of  $n_j$  data points  $(x_{j1}, \dots, x_{jn_j})$ . We assume that the data points in each group are exchangeable, and are to be modeled with a mixture model. While each mixture model has mixing proportions specific to the group, we require that the different groups share the same set of mixture components. The idea is that while different groups have different characteristics given by a different combination of mixing proportions, using the same set of mixture components allows statistical strength to be shared across groups, and allows generalization to new groups.

The HDP is a nonparametric prior which allows the mixture models to share components. It is a distribution over a set of random probability measures over  $(\Theta, \mathcal{B})$ : one probability measure  $G_j$  for each group  $j$ , and a global probability measure  $G_0$ . The global measure  $G_0$  is distributed as  $\text{DP}(\gamma, H)$ , with  $H$  the base measure and  $\gamma$  the concentration parameter, while each  $G_j$  is conditionally independent given  $G_0$ , with distribution  $G_j \sim \text{DP}(\alpha_0, G_0)$ . To complete the description of the HDP mixture model, we associate each  $x_{ji}$  with a factor  $\phi_{ji}$ , with distributions given by  $F(\phi_{ji})$  and  $G_j$  respectively. The overall model is given in Figure 1 left, with conditional distributions:

$$G_0 \mid \gamma, H \sim \text{DP}(\gamma, H) \quad G_j \mid \alpha, G_0 \sim \text{DP}(\alpha_0, G_0) \quad (3)$$

$$\phi_{ji} \mid G_j \sim G_j \quad x_{ji} \mid \phi_{ji} \sim F(\phi_{ji}). \quad (4)$$

The stick-breaking construction (2) shows that  $G_0$  can be expressed as a weighted sum of point masses:  $G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k}$ . This fact that  $G_0$  is atomic plays an important role in ensuring that mixture components are shared across different groups. Since  $G_0$  is the base distribution for the individual  $G_j$ 's, (2) again shows that the atoms of the individual  $G_j$  are samples from  $G_0$ . In particular, since  $G_0$  places non-zero mass only on the atoms  $\theta = (\theta_k)_{k=1}^{\infty}$ , the atoms of  $G_j$  must also come from  $\theta$ , hence we may write:

$$G_0 = \sum_{k=1}^{\infty} \beta_k \delta_{\theta_k} \quad G_j = \sum_{k=1}^{\infty} \pi_{jk} \delta_{\theta_k}. \quad (5)$$

Identifying  $\theta_k$  as the parameters of the  $k^{\text{th}}$  mixture component, we see that each submodel corresponding to distinct groups share the same set of mixture components, but have differing mixing proportions,  $\pi_j = (\pi_{jk})_{k=1}^{\infty}$ .

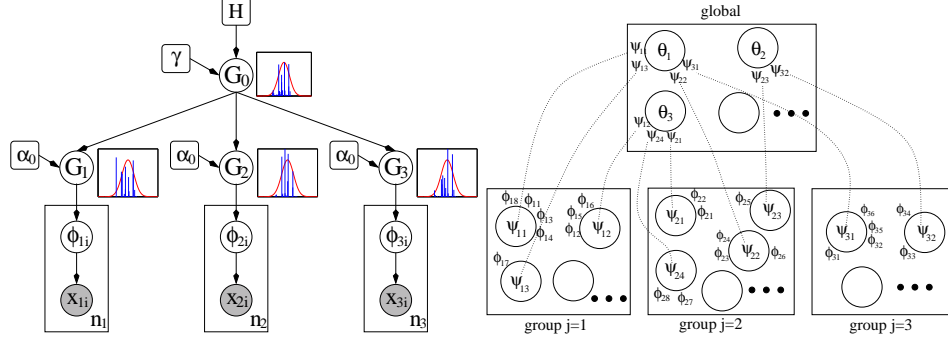


Figure 1: Left: graphical model of an example HDP mixture model with 3 groups. Corresponding to each DP node we also plot a sample draw from the DP using the stick-breaking construction. Right: an instantiation of the CRF representation for the 3 group HDP. Each of the 3 restaurants has customers sitting around tables, and each table is served a dish (which corresponds to customers in the Chinese restaurant for the global DP).

Finally, it is useful to explicitly describe the relationships between the mixing proportions  $\beta$  and  $(\pi_j)_{j=1}^J$ . Details are provided in [10]. Note that the weights  $\pi_j$  are conditionally independent given  $\beta$  since each  $G_j$  is independent given  $G_0$ . Applying (1) to finite partitions of  $\theta$ , we get  $\pi_j \sim \text{DP}(\alpha_0, \beta)$ , where we interpret  $\beta$  and  $\pi_j$  as probability measures over the positive integers. Hence  $\beta$  is simply the putative mixing proportion over the groups. We may in fact obtain an explicit stick-breaking construction for the  $\pi_j$ 's as well. Applying (1) to partitions  $(\{1, \dots, k-1\}, \{k\}, \{k+1, \dots\})$  of positive integers, we have:

$$\pi'_{jk} \sim \text{Beta}(\alpha_0 \beta_k, \alpha_0 (1 - \sum_{l=1}^k \beta_l)) \quad \pi_{jk} = \pi'_{jk} \prod_{l=1}^{k-1} (1 - \pi'_{jl}) . \quad (6)$$

## 4 The Chinese Restaurant Franchise

We describe an alternative view of the HDP based directly upon the distribution a HDP induces on the samples  $\phi_{ji}$ , where we marginalize out  $G_0$  and  $G_j$ 's. This view directly leads to an efficient Gibbs sampler for HDP mixture models, which is detailed in the appendix.

Consider, for one group  $j$ , the distribution of  $\phi_{j1}, \dots, \phi_{jn_j}$  as we marginalize out  $G_j$ . Recall that since  $G_j \sim \text{DP}(\alpha_0, G_0)$  we can describe this distribution by describing how to generate  $\phi_{j1}, \dots, \phi_{jn_j}$  using the CRP. Imagine  $n_j$  customers (each corresponds to a  $\phi_{ji}$ ) at a Chinese restaurant with an unbounded number of tables. The first customer sits at the first table. A subsequent customer sits at an occupied table with probability proportional to the number of customers already there, or at the next unoccupied table with probability proportional to  $\alpha_0$ . Suppose customer  $i$  sat at table  $t_{ji}$ . The conditional distributions are:

$$t_{ji} \mid t_{j1}, \dots, t_{ji-1}, \alpha_0 \sim \sum_t \frac{n_{jt}}{\sum_{t'} n_{jt'} + \alpha_0} \delta_t + \frac{\alpha_0}{\sum_{t'} n_{jt'} + \alpha_0} \delta_{t^{new}} , \quad (7)$$

where  $n_{jt}$  is the number of customers currently at table  $t$ . Once all customers have sat down the seating plan corresponds to a partition of  $\phi_{j1}, \dots, \phi_{jn_j}$ . This is an exchangeable process in that the probability of a partition does not depend on the order in which customers sit down. Now we associate with table  $t$  a draw  $\psi_{jt}$  from  $G_0$ , and assign  $\phi_{ji} = \psi_{jt_{ji}}$ .

Performing this process independently for each group  $j$ , we have now integrated out all the  $G_j$ 's, and have an assignment of each  $\phi_{ji}$  to a sample  $\psi_{jt_{ji}}$  from  $G_0$ , with the partition structures given by CRPs. Notice now that all  $\psi_{jt}$ 's are simply i.i.d. draws from  $G_0$ , which

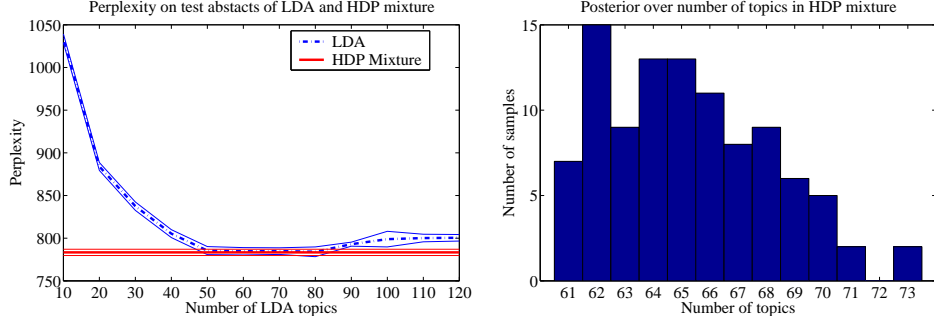


Figure 2: Left: comparison of LDA and HDP mixture. Results are averaged over 10 runs, with error bars being 1 standard error. Right: histogram of the number of topics the HDP mixture used over 100 posterior samples.

is again distributed according to  $DP(\gamma, H)$ , so we may apply the same CRP partitioning process to the  $\psi_{jt}$ 's. Let the customer associated with  $\psi_{jt}$  sit at table  $k_{jt}$ . We have:

$$k_{jt} \mid k_{11}, \dots, k_{1n_1}, k_{21}, \dots, k_{jt-1}, \gamma \sim \sum_k \frac{m_k}{\sum_{k'} m_{jk'} + \gamma} \delta_k + \frac{\gamma}{\sum_{k'} m_{k'} + \alpha_0} \delta_{new} . \quad (8)$$

Finally we associate with table  $k$  a draw  $\theta_k$  from  $H$  and assign  $\psi_{jt} = \theta_{k_{jt}}$ . This completes the generative process for the  $\phi_{ji}$ 's, where we marginalize out  $G_0$  and  $G_j$ 's. We call this generative process the *Chinese restaurant franchise* (CRF). The metaphor is as follows: we have  $J$  restaurants, each with  $n_j$  customers ( $\phi_{ji}$ 's), who sit at tables ( $\psi_{jt}$ 's). Now each table is served a dish ( $\theta_k$ 's) from a menu common to all restaurants. The customers are sociable, preferring large tables with many customers present, and also prefer popular dishes.

## 5 Experiments

We describe 3 experiments in this section to highlight the various aspects of the HDP: its nonparametric nature; its hierarchical nature; and the ease with which we can apply the framework to other models, specifically the HMM.

**Nematode biology abstracts.** To demonstrate the strength of the nonparametric approach as exemplified by the HDP mixture, we compared it against *latent Dirichlet allocation* (LDA), which is a parametric model similar in structure to the HDP [1]. In particular, we applied both models to a corpus of nematode biology abstracts<sup>1</sup>, evaluating the perplexity of both models on held out abstracts. In order to study specifically the nonparametric nature of the HDP, we used the same experimental setup for both models<sup>2</sup>, except that in LDA we had to vary the number of topics used between 10 and 120, while the HDP obtained posterior samples over this automatically.

The results are shown in Figure 2. LDA performs best using between 50 and 80 topics, while the HDP performed just as well as these. Further, the posterior over the number of topics used by HDP is consistent with this range. Notice however that the HDP infers the number of topics automatically, while LDA requires some method of model selection.

<sup>1</sup>Available at <http://elegans.swmed.edu/wli/cgcbib>. There are 5838 abstracts in total. After removing standard stop words and words appearing less than 10 times, we are left with 476441 words in total and a vocabulary size of 5699.

<sup>2</sup>In both models, we used a symmetric Dirichlet distribution with weights of 0.5 for the prior  $H$  over topic distributions, while the concentration parameters are integrated out using a vague gamma prior. Gibbs sampling using the CRF is used.

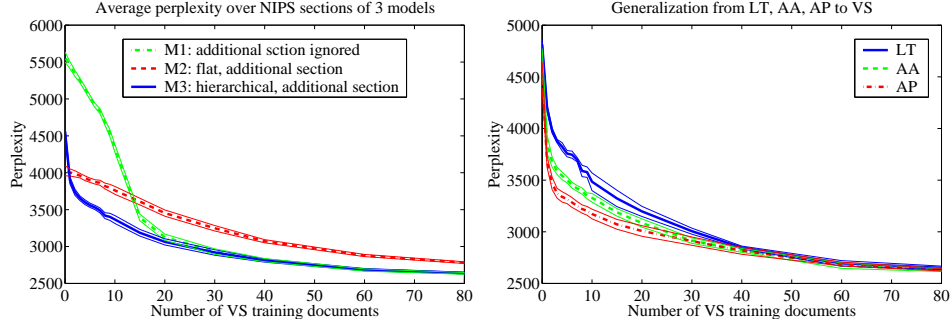


Figure 3: Left: perplexity of test VS documents given training documents from VS and another section for 3 different models. Curves shown are averaged over the other sections and 5 runs. Right: perplexity of test VS documents given LT, AA and AP documents respectively, using M3, averaged over 5 runs. In both, the error bars are 1 standard error.

**NIPS sections.** We applied HDP mixture models to a dataset of NIPS 1-12 papers organized into sections<sup>3</sup>. To highlight the transfer of learning achievable with the HDP, we show improvements to the modeling of a section when the model is also given documents from another section. Our test section is always the VS (vision sciences) section, while the additional section is varied across the other eight. The training set always consist of 80 documents from the other section (so that larger sections like AA (algorithms and architectures) do not get an unfair advantage), plus between 0 and 80 documents from VS. There are 47 test documents, which are held fixed as we vary over the other section and the number  $N$  of training VS documents. We compared 3 different models for this task. The first model (M1) simply ignores documents from the additional section, and uses a HDP to model the VS documents. It serves as a baseline. The second model (M2) uses a HDP mixture model, with one group per document, but lumping together training documents from both sections. The third model (M3) takes a hierarchical approach and models each section separately using a HDP mixture model, and places another DP prior over the common base distributions for both submodels<sup>4</sup>.

As we see in Figure 3 left, the more hierarchical approach of M3 performs best, with perplexity decreasing drastically with modest values of  $N$ , while M1 does worst for small  $N$ . However with increasing  $N$ , M1 improves until it is competitive with M3 but M2 does worst. This is because M2 lumps all the documents together, so is not able to differentiate between the sections, as a result the influence of documents from the other section is unduly strong. This result confirms that the hierarchical approach to the transfer-of-learning problem is a useful one, as it allows useful information to be transferred to a new task (here the modeling of a new section), without the data from the previous tasks overwhelming those in the new task.

We also looked at the performance of the M3 model on VS documents given specific other sections. This is shown in Figure 3 right. As expected, the performance is worst given LT (learning theory), and improves as we move to AA and AP (applications). In Table 1 we show the topics pertinent to VS discovered by the M3 model. First we trained the model

<sup>3</sup>To ensure we are dealing with informative words in the documents, we culled stop words as well as words occurring more than 4000 or less than 50 times in the documents. As sections differ over the years, we assigned by hand the various sections to one of 9 prototypical sections: CS, NS, LT, AA, IM, SP, VS, AP and CN.

<sup>4</sup>Though we have only described the 2 layer HDP the 3 layer extension is straightforward. In fact on our website <http://www.cs.berkeley.edu/~ywtteh/research/npbayes> we have an implementation of the general case where DPs are coupled hierarchically in a tree-structured model.

CS	NS	LT	AA	IM	SP	AP	CN
task representation pattern processing trained representations three process unit patterns	cells cell activity response neuron visual patterns pattern single fig	signal layer gaussian cells fig nonlinearity nonlinear rate eq cell	algorithms test approach methods based point problems form large paper	processing pattern approach architecture single shows simple based large control	visual images video language image pixel acoustic delta lowpass flow	approach based trained test layer features table classification rate paper	ii tree pomdp observable strategy class stochastic history strategies density
examples concept similarity bayesian hypotheses generalization numbers positive classes hypothesis	visual cells cortical orientation receptive contrast spatial cortex stimulus tuning	large examples form point see parameter consider random small optimal	distance tangent image images transformations pattern vectors convolution simard	motion visual velocity flow target chip eye smooth direction optical	signals separation signal sources source matrix blind mixing gradient eq	image images face similarity pixel visual database matching facial examples	policy optimal reinforcement control action states actions step problems goal

Table 1: Topics shared between VS and the other sections. Shown are the two topics with most numbers of VS words, but also with significant numbers of words from the other section.

on all documents from the other section. Then, keeping the assignments of words to topics fixed in the other section, we introduced VS documents and the model decides to reuse some topics from the other section, as well as create new ones. The topics reused by VS documents confirm to our expectations of the overlap between VS and other sections.

**Alice in Wonderland.** The *infinite hidden Markov model* (iHMM) is a nonparametric model for sequential data where the number of hidden states is open-ended and inferred from data [11]. In [10] we show that the HDP framework can be applied to obtain a cleaner formulation of the iHMM<sup>5</sup>, providing effective new inference algorithms and potentially hierarchical extensions. Here we report experimental comparisons of the iHMM against other approaches on sentences taken from Lewis Carroll’s *Alice’s Adventures in Wonderland*.

ML, MAP, and variational Bayesian (VB) [12] models with numbers of states ranging from 1 to 30 were trained multiple times on 20 sentences of average length 51 symbols (27 distinct symbols, consisting of 26 letters and ‘ ’), and tested on 40 sequences of average length 100. Figure 4 shows the perplexity of test sentences. For VB, the predictive probability is intractable to compute, so the modal setting of parameters was used. Both MAP and VB models were given optimal settings of the hyperparameters found in the iHMM. We see that the iHMM has a lower perplexity than every model size for ML, MAP, and VB, and obtains this with one countably infinite model.

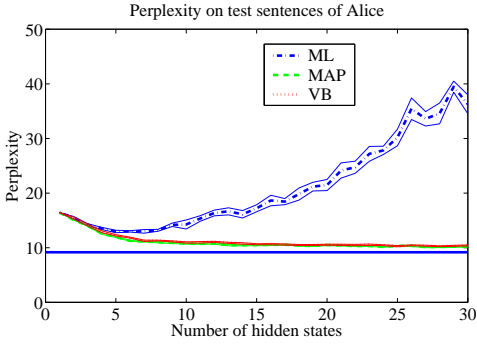


Figure 4: Comparing iHMM (horizontal line) versus ML, MAP and VB trained HMMs. Error bars are 1 standard error (those for iHMM too small to see).

## 6 Discussion

We have described hierarchical Dirichlet processes, a hierarchical, nonparametric model for clustering problems involving multiple groups of data. HDP mixture models are able to automatically determine the appropriate number of mixture components needed, and exhibit sharing of statistical strength across groups by having components shared across groups. We have described the HDP as a distribution over distributions, using both the

<sup>5</sup>In fact the original iHMM paper [11] served as inspiration for this work and first coined the term ‘hierarchical Dirichlet processes’—though their model is not hierarchical in the Bayesian sense, involving priors upon priors, but is rather a set of coupled urn models similar to the CRF.



stick-breaking construction and the Chinese restaurant franchise. In [10] we also describe a fourth perspective based on the infinite limit of finite mixture models, and give detail for how the HDP can be applied to the iHMM. Direct extensions of the model include use of nonparametric priors other than the DP, building higher level hierarchies as in our NIPS experiment, as well as hierarchical extensions to the iHMM.

### Appendix: Gibbs Sampling in the CRF.

The CRF is defined by the variables  $\mathbf{t} = (t_{ji})$ ,  $\mathbf{k} = (k_{jt})$ , and  $\boldsymbol{\theta} = (\theta_k)$ . We describe an inference procedure for the HDP mixture model based on Gibbs sampling  $\mathbf{t}$ ,  $\mathbf{k}$  and  $\boldsymbol{\theta}$  given data items  $\mathbf{x}$ . For the full derivation see [10]. Let  $f(\cdot|\theta)$  and  $h$  be the density functions for  $F(\theta)$  and  $H$  respectively,  $n_{jt}^{-i}$  be the number of  $t_{ji}$ 's equal to  $t$  except  $t_{ji}$ , and  $m_k^{-jt}$  be the number of  $k_{jt}$ 's equal to  $k$  except  $k_{jt}$ . The conditional probability for  $t_{ji}$  given the other variables is proportional to the product of a prior and likelihood term. The prior term is given by (7) where, by exchangeability, we can take  $t_{ji}$  to be the last one assigned. The likelihood is given by  $f(x_{ji}|\theta_{k_{jt}})$  where for  $t = t^{\text{new}}$  we may sample  $k_{jt^{\text{new}}}$  using (8), and  $\theta_{k^{\text{new}}} \sim H$ . The distribution is then:

$$p(t_{ji} = t | \mathbf{t} \setminus t_{ji}, \mathbf{k}, \boldsymbol{\theta}, \mathbf{x}) \propto \begin{cases} \alpha_0 f(x_{ji}|\theta_{k_{jt}}) & \text{if } t = t^{\text{new}} \\ n_{jt}^{-i} f(x_{ji}|\theta_{k_{jt}}) & \text{if } t \text{ currently used.} \end{cases} \quad (9)$$

Similarly the conditional distribution for  $k_{jt}$  is:

$$p(k_{jt} = k | \mathbf{t}, \mathbf{k} \setminus k_{jt}, \boldsymbol{\theta}, \mathbf{x}) \propto \begin{cases} \gamma \prod_{i:t_{ji}=t} f(x_{ji}|\theta_k) & \text{if } k = k^{\text{new}} \\ m_k^{-t} \prod_{i:t_{ji}=t} f(x_{ji}|\theta_k) & \text{if } k \text{ currently used.} \end{cases} \quad (10)$$

where  $\theta_{k^{\text{new}}} \sim H$ . Finally the conditional distribution for  $\theta_k$  is:

$$p(\theta_k | \mathbf{t}, \mathbf{k}, \boldsymbol{\theta} \setminus \theta_k, \mathbf{x}) \propto h(\theta_k) \prod_{j:i:k_{jt_{ji}}=k} f(x_{ji}|\theta_k) \quad (11)$$

If  $H$  is conjugate to  $F(\cdot)$  we have the option of integrating out  $\boldsymbol{\theta}$ .

### References

- [1] D.M. Blei, A.Y. Ng, and M.I. Jordan. Latent Dirichlet allocation. *JMLR*, 3:993–1022, 2003.
- [2] M.D. Escobar and M. West. Bayesian density estimation and inference using mixtures. *Journal of the American Statistical Association*, 90:577–588, 1995.
- [3] S.N. MacEachern and P. Müller. Estimating mixture of Dirichlet process models. *Journal of Computational and Graphical Statistics*, 7:223–238, 1998.
- [4] T.S. Ferguson. A Bayesian analysis of some nonparametric problems. *Annals of Statistics*, 1(2):209–230, 1973.
- [5] D. Aldous. Exchangeability and related topics. In *École d'été de probabilités de Saint-Flour XIII–1983*, pages 1–198. Springer, Berlin, 1985.
- [6] J. Sethuraman. A constructive definition of Dirichlet priors. *Statistica Sinica*, 4:639–650, 1994.
- [7] R.M. Neal. Markov chain sampling methods for Dirichlet process mixture models. *Journal of Computational and Graphical Statistics*, 9:249–265, 2000.
- [8] C.E. Rasmussen. The infinite Gaussian mixture model. In *NIPS*, volume 12, 2000.
- [9] D.M. Blei, T.L. Griffiths, M.I. Jordan, and J.B. Tenenbaum. Hierarchical topic models and the nested Chinese restaurant process. *NIPS*, 2004.
- [10] Y.W. Teh, M.I. Jordan, M.J. Beal, and D.M. Blei. Hierarchical dirichlet processes. Technical Report 653, Department of Statistics, University of California at Berkeley, 2004.
- [11] M.J. Beal, Z. Ghahramani, and C.E. Rasmussen. The infinite hidden Markov model. In *NIPS*, volume 14, 2002.
- [12] M.J. Beal. *Variational Algorithms for Approximate Bayesian Inference*. PhD thesis, Gatsby Unit, University College London, 2004.

# HIERARCHICAL DENSITY ORDER EMBEDDINGS

**Ben Athiwaratkun, Andrew Gordon Wilson**

Cornell University  
Ithaca, NY 14850, USA

## ABSTRACT

By representing words with probability densities rather than point vectors, probabilistic word embeddings can capture rich and interpretable semantic information and uncertainty. The uncertainty information can be particularly meaningful in capturing *entailment* relationships – whereby general words such as “entity” correspond to broad distributions that encompass more specific words such as “animal” or “instrument”. We introduce *density order embeddings*, which learn hierarchical representations through encapsulation of probability densities. In particular, we propose simple yet effective loss functions and distance metrics, as well as graph-based schemes to select negative samples to better learn hierarchical density representations. Our approach provides state-of-the-art performance on the WORD-NET hypernym relationship prediction task and the challenging HYPERLEX lexical entailment dataset – while retaining a rich and interpretable density representation.

## 1 INTRODUCTION

Learning feature representations of natural data such as text and images has become increasingly important for understanding real-world concepts. These representations are useful for many tasks, ranging from semantic understanding of words and sentences (Mikolov et al., 2013; Kiros et al., 2015), image caption generation (Vinyals et al., 2015), textual entailment prediction (Rocktäschel et al., 2015), to language communication with robots (Bisk et al., 2016).

Meaningful representations of text and images capture visual-semantic information, such as hierarchical structure where certain entities are abstractions of others. For instance, an image caption “A dog and a frisbee” is an abstraction of many images with possible lower-level details such as a dog jumping to catch a frisbee or a dog sitting with a frisbee (Figure 1a). A general word such as “object” is also an abstraction of more specific words such as “house” or “pool”. Recent work by Vendrov et al. (2016) proposes learning such asymmetric relationships with *order embeddings* – vector representations of non-negative coordinates with partial order structure. These embeddings are shown to be effective for word hypernym classification, image-caption ranking and textual entailment (Vendrov et al., 2016).

Another recent line of work uses probability distributions as rich feature representations that can capture the semantics and uncertainties of concepts, such as Gaussian word embeddings (Vilnis & McCallum, 2015), or extract multiple meanings via multimodal densities (Athiwaratkun & Wilson, 2017). Probability distributions are also natural at capturing orders and are suitable for tasks that involve hierarchical structures. An abstract entity such as “animal” that can represent specific entities such as “insect”, “dog”, “bird” corresponds to a broad distribution, encapsulating the distributions for these specific entities. For example, in Figure 1c, the distribution for “insect” is more concentrated than for “animal”, with a high density occupying a small volume in space.

Such entailment patterns can be observed from density word embeddings through *unsupervised* training based on word contexts (Vilnis & McCallum, 2015; Athiwaratkun & Wilson, 2017). In the unsupervised settings, density embeddings are learned via maximizing the similarity scores between nearby words. In these cases, the density encapsulation behavior arises due to the word occurrence pattern that a general word can often substitute more specific words; for instance, the word “tea” in a sentence “I like iced tea” can be substituted by “beverages”, yielding another natural sentence “I like iced beverages”. Therefore, the probability density of a general concept such as “beverages” tends to have a larger variance than specific ones such as “tea”, reflecting higher uncertainty in meanings

since a general word can be used in many contexts. However, the information from word occurrences alone is not sufficient to train meaningful embeddings of some concepts. For instance, it is fairly common to observe sentences “Look at the cat”, or “Look at the dog”, but not “Look at the mammal”. Therefore, due to the way we typically express natural language, it is unlikely that the word “mammal” would be learned as a distribution that encompasses both “cat” and “dog”, since “mammal” rarely occurs in similar contexts.

Rather than relying on the information from word occurrences, one can do *supervised* training of density embeddings on hierarchical data. In this paper, we propose new training methodology to enable effective supervised probabilistic density embeddings. Despite providing rich and intuitive word representations, with a natural ability to represent order relationships, probabilistic embeddings have only been considered in a small number of pioneering works such as Vilnis & McCallum (2015), and these works are almost exclusively focused on *unsupervised embeddings*. Probabilistic Gaussian embeddings trained directly on labeled data have been briefly considered but perform surprisingly poorly compared to other competing models (Vendrov et al., 2016; Vulić et al., 2016).

Our work reaches a very different conclusion: probabilistic Gaussian embeddings can be *highly effective* at capturing ordering and are suitable for modeling hierarchical structures, and can even achieve state-of-the-art results on hypernym prediction and graded lexical entailment tasks, so long as one uses the right training procedures.

In particular, we make the following contributions.

- (a) We adopt a new form of loss function for training hierarchical probabilistic order embeddings.
- (b) We introduce the notion of soft probabilistic encapsulation orders and a thresholded divergence-based penalty function, which do not over-penalize words with a sufficient encapsulation.
- (c) We introduce a new graph-based scheme to select negative samples to contrast the true relationship pairs during training. This approach incorporates hierarchy information to the negative samples that help facilitate training and has added benefits over the hierarchy-agnostic sampling schemes previously used in literature.
- (d) We also demonstrate that initializing the right variance scale is highly important for modeling hierarchical data via distributions, allowing the model to exhibit meaningful encapsulation orders.

The outline of our paper is as follows. In Section 2, we introduce the background for Gaussian embeddings (Vilnis & McCallum, 2015) and vector order embeddings (Vendrov et al., 2016). We describe our training methodology in Section 3, where we introduce the notion of soft encapsulation orders (Section 3.2) and explore different divergence measures such as the expected likelihood kernel, KL divergence, and a family of Rényi alpha divergences (Section 3.3). We describe the experiment details in Section 4 and offer a qualitative evaluation of the model in Section 4.3, where we show the visualization of the density encapsulation behavior. We show quantitative results on the WORDNET Hypernym prediction task in Section 4.2 and a graded entailment dataset HYPERLEX in Section 4.4.

In addition, we conduct experiments to show that our proposed changes to learn Gaussian embeddings contribute to the increased performance. We demonstrate (a) the effects of our loss function in Section A.2.3, (b) soft encapsulation in Section A.2.1, (c) negative sample selection in Section 4.4], and (d) initial variance scale in Section A.2.2.

We make our code publicly available.<sup>1</sup>

## 2 BACKGROUND AND RELATED WORK

### 2.1 GAUSSIAN EMBEDDINGS

Vilnis & McCallum (2015) was the first to propose using probability densities as word embeddings. In particular, each word is modeled as a Gaussian distribution, where the mean vector represents the semantics and the covariance describes the uncertainty or nuances in the meanings. These embeddings are trained on a natural text corpus by maximizing the similarity between words that are in the same local context of sentences. Given a word  $w$  with a true context word  $c_p$  and a randomly sampled

<sup>1</sup><https://github.com/benathi/density-order-emb>

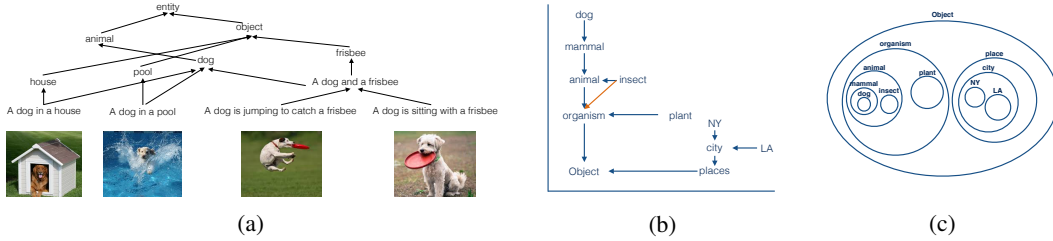


Figure 1: (a) Captions and images in the visual-semantic hierarchy. (b) Vector order embedding (Vendrov et al., 2016) where specific entities have higher coordinate values. (c) Density order embedding where specific entities correspond to concentrated distributions encapsulated in broader distributions of general entities.

word  $c_n$  (negative context), Gaussian embeddings are learned by minimizing the rank objective in Equation 1, which pushes the similarity of the true context pair  $E(w, c_p)$  above that of the negative context pair  $E(w, c_n)$  by a margin  $m$ .

$$L_m(w, c_p, c_n) = \max(0, m - E(w, c_p) + E(w, c_n)) \quad (1)$$

The similarity score  $E(u, v)$  for words  $u, v$  can be either  $E(u, v) = -\text{KL}(f_u, f_v)$  or  $E(u, v) = \log \langle f_u, f_v \rangle_{L_2}$  where  $f_u, f_v$  are the distributions of words  $u$  and  $v$ , respectively. The Gaussian word embeddings contain rich semantic information and performs competitively in many word similarity benchmarks.

The true context word pairs  $(w, c_p)$  are obtained from natural sentences in a text corpus such as Wikipedia. In some cases, specific words can be replaced by a general word in a similar context. For instance, “I love cats” or “I love dogs” can be replaced with “I love animals”. Therefore, the trained word embeddings exhibit lexical entailment patterns where specific words such as “dog” and “cat” are concentrated distributions that are encompassed by a more dispersed distribution of “animal”, a word that “cat” and “dog” entail. The broad distribution of a general word agrees with the *distributional informativeness hypothesis* proposed by Santus et al. (2014), which says that a generic word can occur in more general contexts in place of the specific ones that entail it.

However, some word entailment pairs have weak density encapsulation patterns due to the nature of word diction. For instance, even though “dog” and “cat” both entail “mammal”, it is rarely the case that we observe a sentence “I have a mammal” as opposed to “I have a cat” in a natural corpus; therefore, after training density word embeddings on word occurrences, encapsulation of some true entailment instances do not occur.

## 2.2 PARTIAL ORDERS AND VECTOR ORDER EMBEDDINGS

We describe the concepts of partial orders and vector order embeddings proposed by Vendrov et al. (2016), which we will later consider in the context of our hierarchical density order embeddings.

A partial order over a set of points  $X$  is a binary relation  $\preceq$  such that for  $a, b, c \in X$ , the following properties hold: (1)  $a \preceq a$  (reflexivity); (2) if  $a \preceq b$  and  $b \preceq a$  then  $a = b$  (antisymmetry); and (3) if  $a \preceq b$  and  $b \preceq c$  then  $a \preceq c$  (transitivity). An example of a partially ordered set is a set of nodes in a tree where  $a \preceq b$  means  $a$  is a child node of  $b$ . This concept has applications in natural data such as lexical entailment. For words  $a$  and  $b$ ,  $a \preceq b$  means that every instance of  $a$  is an instance of  $b$ , or we can say that  $a$  entails  $b$ . We also say that  $(a, b)$  has a *hypernym* relationship where  $a$  is a hyponym of  $b$  and  $b$  is a hypernym of  $a$ . This relationship is asymmetric since  $a \preceq b$  does not necessarily imply  $b \preceq a$ . For instance,  $\text{aircraft} \preceq \text{vehicle}$  but it is not true that  $\text{vehicle} \preceq \text{aircraft}$ .

An order-embedding is a function  $f : (X, \preceq_X) \rightarrow (Y, \preceq_Y)$  where  $a \preceq_X b$  if and only if  $f(a) \preceq_Y f(b)$ . Vendrov et al. (2016) proposes to learn the embedding  $f$  on  $Y = \mathbb{R}_+^N$  where all coordinates are non-negative. Under  $\mathbb{R}_+^N$ , there exists a partial order relation called the *reversed product order on  $\mathbb{R}_+^N$* :  $x \preceq y$  if and only if  $\forall i, x_i \geq y_i$ . That is, a point  $x$  entails  $y$  if and only if all the coordinate values of  $x$  is higher than  $y$ ’s. The origin represents the most general entity at the top of the order hierarchy and the points further away from the origin become more specific. Figure 1b demonstrates the vector order embeddings on  $\mathbb{R}_+^N$ . We can see that since  $\text{insect} \preceq \text{animal}$  and  $\text{animal} \preceq \text{organism}$ , we can

infer directly from the embedding that  $\text{insect} \preceq \text{organism}$  (orange line, diagonal line). To learn the embeddings, Vendrov et al. (2016) proposes a penalty function  $E(x, y) = \|\max(0, y - x)\|^2$  for a pair  $x \preceq y$  which has the property that it is positive if and only if the order is violated.

### 2.3 OTHER RELATED WORK

Li et al. (2017) extends Vendrov et al. (2016) for knowledge representation on data such as ConceptNet (Speer et al., 2016). Another related work by Hockenmaier & Lai (2017) embeds words and phrases in a vector space and uses denotational probabilities for textual entailment tasks. Our models offer an improvement on order embeddings and can be applicable to such tasks, which view as a promising direction for future work.

## 3 METHODOLOGY

In Section 3.1, we describe the partial orders that can be induced by density encapsulation. Section 3.2 describes our training approach that softens the notion of strict encapsulation with a viable penalty function.

### 3.1 STRICT ENCAPSULATION PARTIAL ORDERS

A partial order on probability densities can be obtained by the notion of encapsulation. That is, a density  $f$  is more specific than a density  $g$  if  $f$  is encompassed in  $g$ . The degree of encapsulation can vary, which gives rise to multiple order relations. We define an order relation  $\preceq_\eta$  for  $\eta \geq 0$  where  $\eta$  indicates the degree of encapsulation required for one distribution to entail another. More precisely, for distributions  $f$  and  $g$ ,

$$f \preceq_\eta g \Leftrightarrow \{x : f(x) > \eta\} \subseteq \{x : g(x) > \eta\}. \quad (2)$$

Note that  $\{x : f(x) > \eta\}$  is a set where the density  $f$  is greater than the threshold  $\eta$ . The relation in Equation 2 says that  $f$  entails  $g$  if and only if the set of  $g$  contains that of  $f$ . In Figure 2, we depict two Gaussian distributions with different mean vectors and covariance matrices. Figure 2 (left) shows the density values of distributions  $f$  (narrow, blue) and  $g$  (broad, orange) and different threshold levels. Figure 2 (right) shows that different  $\eta$ 's give rise to different partial orders. For instance, we observe that neither  $f \preceq_{\eta_1} g$  nor  $g \preceq_{\eta_1} f$  but  $f \preceq_{\eta_3} g$ .

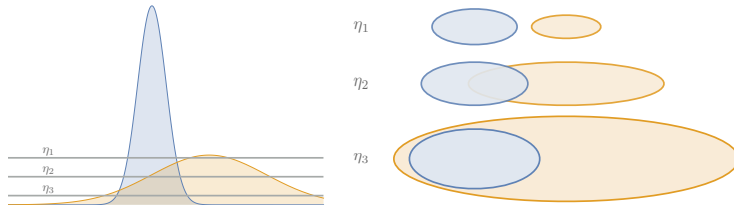


Figure 2: Strict encapsulation orders induced by different  $\eta$  values.

### 3.2 SOFT ENCAPSULATION ORDERS

A plausible penalty function for the order relation  $f \preceq_\eta g$  is a set measure on  $\{x : f(x) > \eta\} - \{x : g(x) > \eta\}$ . However, this penalty is difficult to evaluate for most distributions, including Gaussians. Instead, we use simple penalty functions based on asymmetric divergence measures between probability densities. Divergence measures  $D(\cdot||\cdot)$  have a property that  $D(f||g) = 0$  if and only if  $f = g$ . Using  $D(\cdot||\cdot)$  to represent order violation is undesirable since the penalty should be 0 if  $f \neq g$  but  $f \preceq g$ . Therefore, we propose using a thresholded divergence

$$d_\gamma(f, g) = \max(0, D(f||g) - \gamma),$$

which can be zero if  $f$  is properly encapsulated in  $g$ . We discuss the effectiveness of using divergence thresholds in Section A.2.1.



We note that by using  $d_\gamma(\cdot, \cdot)$  as a violation penalty, we no longer have the strict *partial order*. In particular, the notion of transitivity in a partial order is not guaranteed. For instance, if  $f \preceq g$  and  $g \preceq h$ , our density order embeddings would yield  $d_\gamma(f, g) = 0$  and  $d_\gamma(g, h) = 0$ . However, it is not necessarily the case that  $d_\gamma(f, h) = 0$  since  $D(f||h)$  can be greater than  $\gamma$ . This is not a drawback since a high value of  $D(f||h)$  reflects that the hypernym relationship is not direct, requiring many edges from  $f$  to  $h$  in the hierarchy. The extent of encapsulation contains useful entailment information, as demonstrated in Section 4.4 where our model scores highly correlate with the annotated scores of a challenging lexical entailment dataset and achieves state-of-the-art results.

Another property, antisymmetry, does not strictly hold since if  $d_\gamma(f, g) = 0$  and  $d_\gamma(g, f) = 0$  does not imply  $f = g$ . However, in this situation, it is necessary that  $f$  and  $g$  overlap significantly if  $\gamma$  is small. Due to the fact that the  $d_\gamma(\cdot, \cdot)$  does not strictly induce a partial order, we refer to this model as *soft density order embeddings* or simply *density order embeddings*.

### 3.3 DIVERGENCE MEASURES

#### 3.3.1 ASYMMETRIC DIVERGENCE

**Kullback-Leibler (KL) Divergence** The KL divergence is an asymmetric measure of the difference between probability distributions. For distributions  $f$  and  $g$ ,  $\text{KL}(g||f) \equiv \int g(x) \log \frac{g(x)}{f(x)} dx$  imposes a high penalty when there is a region of points  $x$  such that the density  $f(x)$  is low but  $g(x)$  is high. An example of such a region is the area on the left of  $f$  in Figure 2. This measure penalizes the situation where  $f$  is a concentrated distribution relative to  $g$ ; that is, if the distribution  $f$  is encompassed by  $g$ , then the KL yields high penalty. For  $d$ -dimensional Gaussians  $f = \mathcal{N}_d(\mu_f, \Sigma_f)$  and  $g = \mathcal{N}_d(\mu_g, \Sigma_g)$ ,

$$2D_{KL}(f||g) = \log(\det(\Sigma_g)/\det(\Sigma_f)) - d + \text{tr}(\Sigma_g^{-1}\Sigma_f) + (\mu_f - \mu_g)^T \Sigma_g^{-1}(\mu_f - \mu_g) \quad (3)$$

**Rényi  $\alpha$ -Divergence** is a general family of divergence with varying scale of zero-forcing penalty (Rényi, 1961). Equation 4 describes the general form of the  $\alpha$ -divergence for  $\alpha \neq 0, 1$  (Liese & Vajda, 1987). We note that for  $\alpha \rightarrow 0$  or 1, we recover the KL divergence and the reverse KL divergence; that is,  $\lim_{\alpha \rightarrow 1} D_\alpha(f||g) = \text{KL}(f||g)$  and  $\lim_{\alpha \rightarrow 0} D_\alpha(f||g) = \text{KL}(g||f)$  (Pardo, 2006). The  $\alpha$ -divergences are asymmetric for all  $\alpha$ 's, except for  $\alpha = 1/2$ .

$$D_\alpha(f||g) = \frac{1}{\alpha(\alpha-1)} \log \left( \int \frac{f(x)^\alpha}{g(x)^{\alpha-1}} dx \right) \quad (4)$$

For two multivariate Gaussians  $f$  and  $g$ , we can write the Rényi divergence as (Pardo, 2006):

$$2D_\alpha(f||g) = -\frac{1}{\alpha(\alpha-1)} \log \frac{\det(\alpha\Sigma_g + (1-\alpha)\Sigma_f)}{(\det(\Sigma_f)^{1-\alpha} \cdot \det(\Sigma_g)^\alpha)} + (\mu_f - \mu_g)^T (\alpha\Sigma_g + (1-\alpha)\Sigma_f)^{-1} (\mu_f - \mu_g). \quad (5)$$

The parameter  $\alpha$  controls the degree of *zero forcing* where minimizing  $D_\alpha(f||g)$  for high  $\alpha$  results in  $f$  being more concentrated to the region of  $g$  with high density. For low  $\alpha$ ,  $f$  tends to be *mass-covering*, encompassing regions of  $g$  including the low density regions. Recent work by Li & Turner (2016) demonstrates that different applications can require different degrees of zero-forcing penalty.

#### 3.3.2 SYMMETRIC DIVERGENCE

**Expected Likelihood Kernel** The expected likelihood kernel (ELK) (Jebara et al., 2004) is a symmetric measure of affinity, define as  $K(f, g) = \langle f, g \rangle_{\mathcal{H}}$ . For two Gaussians  $f$  and  $g$ ,

$$2 \log \langle f, g \rangle_{\mathcal{H}} = -\log \det(\Sigma_f + \Sigma_g) - d \log(2\pi) - (\mu_f - \mu_g)^T (\Sigma_f + \Sigma_g)^{-1} (\mu_f - \mu_g) \quad (6)$$

Since this kernel is a similarity score, we use its negative as our penalty. That is,  $D_{\text{ELK}}(f||g) = -2 \log \langle f, g \rangle_{\mathcal{H}}$ . Intuitively, the asymmetric measures should be more successful at training density order embeddings. However, a symmetric measure can result in a correct encapsulation order as well, since a general entity often has to minimize the penalty with many specific elements and consequently ends up having a broad distribution to lower the average loss. The expected likelihood kernel is used to train Gaussian and Gaussian Mixture word embeddings on a large text corpus (Vilnis & McCallum, 2015; Athiwaratkun & Wilson, 2017) where the model performs well on the word entailment dataset (Baroni et al., 2012).

### 3.4 LOSS FUNCTION

To learn our density embeddings, we use a loss function similar to that of Vendrov et al. (2016). Minimizing this function (Equation 7) is equivalent to minimizing the penalty between a true relationship pair  $(u, v)$  where  $u \preceq v$ , but pushing the penalty to be above a margin  $m$  for the negative example  $(u', v')$  where  $u' \not\preceq v'$ :

$$\sum_{(u,v) \in \mathcal{D}} d(u, v) + \max\{0, m - d(u', v')\} \quad (7)$$

We note that this loss function is different than the rank-margin loss introduced in the original Gaussian embeddings (Equation 1). Equation 7 aims to reduce the dissimilarity of a true relationship pair  $d(u, v)$  with no constraint, unlike in Equation 1, which becomes zero if  $d(u, v)$  is above  $d(u', v')$  by margin  $m$ .

### 3.5 SELECTING NEGATIVE SAMPLES

In many embedding models such as WORD2VEC (Mikolov et al., 2013) or Gaussian embeddings (Vilnis & McCallum, 2015), negative samples are often used in the training procedure to contrast with true samples from the dataset. For flat data such as words in a text corpus, negative samples are selected randomly from a unigram distribution. We propose new graph-based methods to select negative samples that are suitable for hierarchical data, as demonstrated by the improved performance of our density embeddings. In our experiments, we use various combinations of the following methods.

**Method S1:** A simple negative sampling procedure used by Vendrov et al. (2016) is to replace a true hypernym pair  $(u, v)$  with either  $(u, v')$  or  $(u', v)$  where  $u', v'$  are randomly sampled from a uniform distribution of vertices. **Method S2:** We use a negative sample  $(v, u)$  if  $(u, v)$  is a true relationship pair, to make  $D(v||u)$  higher than  $D(u||v)$  in order to distinguish the directionality of density encapsulation. **Method S3:** It is important to increase the divergence between neighbor entities that do not entail each other. Let  $A(w)$  denote all descendants of  $w$  in the training set  $\mathcal{D}$ , including  $w$  itself. We first randomly sample an entity  $w \in \mathcal{D}$  that has at least 2 descendants and randomly select a descendant  $u \in A(w) - \{w\}$ . Then, we randomly select an entity  $v \in A(w) - A(u)$  and use the random neighbor pair  $(v, u)$  as a negative sample. Note that we can have  $u \preceq v$ , in which case the pair  $(v, u)$  is a reverse relationship. **Method S4:** Same as S3 except that we sample  $v \in A(w) - A(u) - \{w\}$ , which excludes the possibility of drawing  $(w, u)$ .

## 4 EXPERIMENTS

We have introduced density order embeddings (DOE) to model hierarchical data via encapsulation of probability densities. We propose using a new loss function, graph-based negative sample selections, and a penalty relaxation to induce soft partial orders. In this section, we show the effectiveness of our model on WORDNET hypernym prediction and a challenging graded lexical entailment task, where we achieve state-of-the-art performance.

First, we provide the training details in Section 4.1 and describe the hypernym prediction experiment in 4.2. We offers insights into our model with the qualitative analysis and visualization in Section 4.3. We evaluate our model on HYPERLEX, a lexical entailment dataset in Section 4.4.

### 4.1 TRAINING DETAILS

We have a similar data setup to the experiment by Vendrov et al. (2016) where we use the transitive closure of WORDNET noun hypernym relationships which contains 82, 115 synsets and 837, 888 hypernym pairs from 84, 427 direct hypernym edges. We obtain the data using the WORDNET API of NLTK version 3.2.1 (Loper & Bird, 2002).

The validation set contains 4000 true hypernym relationships as well as 4000 false hypernym relationships where the false hypernym relationships are constructed from the S1 negative sampling described in Section 3.5. The same process applies for the test set with another set of 4000 true hypernym relationships and 4000 false hypernym relationships.

We use  $d$ -dimensional Gaussian distributions with diagonal covariance matrices. We use  $d = 50$  as the default dimension and analyze the results using different  $d$ 's in Section A.2.4. We initialize the mean vectors to have a unit norm and normalize the mean vectors in the training graph. We initialize the diagonal variance components to be all equal to  $\beta$  and optimize on the unconstrained space of  $\log(\Sigma)$ . We discuss the important effects of the initial variance scale in Section A.2.2.

We use a minibatch size of 500 true hypernym pairs and use varying number of negative hypernym pairs, depending on the negative sample combination proposed in Section 3.5. We discuss the results for many selection strategies in Section 4.4. We also experiment with multiple divergence measures  $D(\cdot||\cdot)$  described in Section 3.3. We use  $D(\cdot||\cdot) = D_{KL}(\cdot||\cdot)$  unless stated otherwise. Section A.2.5 considers the results using the  $\alpha$ -divergence family with varying degrees of zero-forcing parameter  $\alpha$ 's. We use the Adam optimizer (Kingma & Ba, 2014) and train our model for at most 20 epochs. For each energy function, we tune the hyperparameters on grids. The hyperparameters are the loss margin  $m$ , the initial variance scale  $\beta$ , and the energy threshold  $\gamma$ . We evaluate the results by computing the penalty on the validation set to find the best threshold for binary classification, and use this threshold to perform prediction on the test set. Section A.1 describes the hyperparameters for all our models.

## 4.2 HYPERNYM PREDICTION

We show the prediction accuracy results on the test set of WORDNET hypernyms in Table 1. We compare our results with **vector order-embeddings** (VOE) by Vendrov et al. (2016) (VOE model details are in Section 2.2). Another important baseline is the **transitive closure**, which requires no learning and classifies if a held-out edge is a hypernym relationship by determining if it is in the union of the training edges. **word2gauss** and **word2gauss<sup>†</sup>** are the Gaussian embeddings trained using the loss function in Vilnis & McCallum (2015) (Equation 1) where **word2gauss** is the result reported by Vendrov et al. (2016) and **word2gauss<sup>†</sup>** is the best performance of our replication (see Section A.2.3 for more details). Our density order embedding (DOE) outperforms the implementation by Vilnis & McCallum (2015) significantly; this result highlights the fact that our different approach for training Gaussian embeddings can be crucial to learning hierarchical representations.

We observe that the symmetric model (ELK) performs quite well for this task despite the fact that the symmetric metric cannot capture directionality. In particular, ELK can accurately detect pairs of concepts with no relationships when they're far away in the density space. In addition, for pairs that are related, ELK can detect pairs that overlap significantly in density space. The lack of directionality has more pronounced effects in the graded lexical entailment task (Section 4.4) where we observe a high degradation in performance if ELK is used instead of KL.

We find that our method outperforms vector order embeddings (VOE) (Vendrov et al., 2016). We also find that DOE is very strong in a 2-dimensional Gaussian embedding example, trained for the purpose of visualization in Section 4.3, despite only having only 4 parameters: 2 from 2-dimensional  $\mu$  and another 2 from the diagonal  $\Sigma$ . The results of DOE using a symmetric measure also outperforms the baselines on this experiment, but has a slightly lower accuracy than the asymmetric model.

Figure 3 offers an explanation as to why our density order embeddings might be easier to learn, compared to the vector counterpart. In certain cases such as fitting a general concept `entity` to the embedding space, we simply need to adjust the distribution of `entity` to be broad enough to encompass all other concepts. In the vector counterpart, it might be required to shift many points further from the origin to accommodate `entity` to reduce cascading order violations.



Figure 3: **(Left)** Adding a concept `entity` to vector order embedding **(Right)** Adding a concept `entity` to density order embedding

Table 1: Classification accuracy on hypernym relationship test set from WordNet.

Method	Test Accuracy (%)
transitive closure	88.2
word2gauss	86.6
word2gauss†	88.6
VOE (symmetric)	84.2
VOE	90.6
DOE (ELK)	92.1
DOE (KL, reversed)	83.2
DOE (KL)	<b>92.3</b>
DOE (KL, $d = 2$ )	89.2

### 4.3 QUALITATIVE ANALYSIS

For qualitative analysis, we additionally train a 2-dimensional Gaussian model for visualization. Our qualitative analysis shows that the encapsulation behavior can be observed in the trained model. Figure 4 demonstrates the ordering of synsets in the density space. Each ellipse represents a Gaussian distribution where the center is given by the mean vector  $\mu$  and the major and minor axes are given by the diagonal standard deviations  $\sqrt{\Sigma}$ , scaled by 300 for the  $x$  axis and 30 for the  $y$  axis, for visibility.

Most hypernym relationships exhibit encapsulation behavior where the hypernym encompasses the synset that entails it. For instance, the distribution of `whole.n.02` is subsumed in the distribution of `physical_entity.n.01`. Note that `location.n.01` is not entirely encapsulated by `physical_entity.n.01` under this visualization. However, we can still predict which entity should be the hypernym among the two since the KL divergence of one given another would be drastically different. This is because a large part of `physical_entity.n.01` has considerable density at the locations where `location.n.01` has very low density. This causes  $\text{KL}(\text{physical\_entity.n.01} \parallel \text{location.n.01})$  to be very high (5103) relative to  $\text{KL}(\text{location.n.01} \parallel \text{physical\_entity.n.01})$  (206). Table 2 shows the KL values for all pairs where we note that the numbers are from the full model ( $d = 50$ ).

Another interesting pair is `city.n.01`  $\preceq$  `location.n.01` where we see the two distributions have very similar contours and the encapsulation is not as distinct. In our full model  $d = 50$ , the distribution of `location.n.01` encompasses `city.n.01`’s, indicated by low  $\text{KL}(\text{city.n.01} \parallel \text{location.n.01})$  but high  $\text{KL}(\text{location.n.01} \parallel \text{city.n.01})$ .

Figure 4 (Right) demonstrates the idea that synsets on the top of the hypernym hierarchy usually have higher “volume”. A convenient metric that reflects this quantity is  $\log \det(\Sigma)$  for a Gaussian distribution with covariance  $\Sigma$ . We can see that the synset, `physical_entity.n.01`, being the hypernym of all the synsets shown, has the highest  $\log \det(\Sigma)$  whereas entities that are more specific such as `object.n.01`, `whole.n.02` and `living_thing` have decreasingly lower volume.

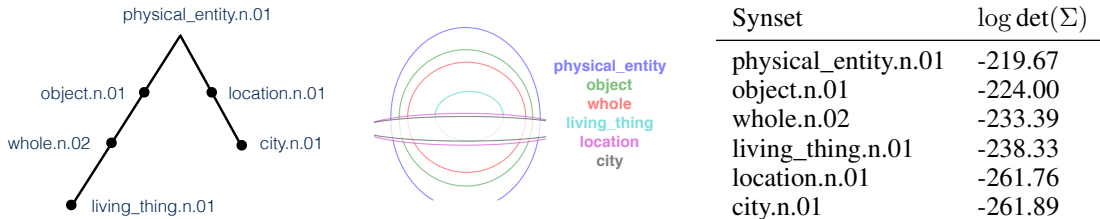


Figure 4: [best viewed electronically] (Left) Synsets and their hypernym relationships from WordNet. (Middle) Visualization of our 2-dimensional Gaussian order embedding. (Right) The Gaussian “volume” ( $\log \det \Sigma$ ) of the 50-dimensional Gaussian model.

Table 2:  $KL(\text{column}||\text{row})$ . Cells in boldface indicate true WORDNET hypernym relationships ( $\text{column} \preceq \text{row}$ ). Our model predicts a synset pair as a hypernym if the KL less than 1900, where this value is tuned based on the validation set. Most relationship pairs are correctly predicted except for the underlined cells.

	city	location	living_thing	whole	object	physical_entity
city	0	<u>1025</u>	4999	4673	23673	4639
location	<b>159</b>	0	4324	4122	26121	5103
living_thing	3623	6798	0	<u>1452</u>	2953	5936
whole	3033	6367	<b>66</b>	0	6439	6682
object	<u>138</u>	<u>80</u>	<b>125</b>	<b>77</b>	0	6618
physical_entity	<b>232</b>	<b>206</b>	<b>193</b>	<b>166</b>	<b>152</b>	0

#### 4.4 GRADED LEXICAL ENTAILMENT

HYPERLEX is a lexical entailment dataset which has fine-grained human annotated scores between concept pairs, capturing varying degrees of entailment (Vulić et al., 2016). Concept pairs in HYPERLEX reflect many variants of hypernym relationships, such as `no-rel` (no lexical relationship), `ant` (antonyms), `syn` (synonyms), `cohyp` (sharing a hypernym but not a hypernym of each other), `hyp` (hypernym), `rhyp` (reverse hypernym). We use the noun dataset of HYPERLEX for evaluation, which contains 2,163 pairs.

We evaluate our model by comparing our model scores against the annotated scores. Obtaining a high correlation on a fine-grained annotated dataset is a much harder task compared to a binary prediction, since performing well requires meaningful model scores in order to reflect nuances in hypernymy. We use negative divergence as our score for hypernymy scale where large values indicate high degrees of entailment.

We note that the concepts in our trained models are WORDNET synsets, where each synset corresponds to a specific meaning of a word. For instance, `pop.n.03` has a definition “a sharp explosive sound as from a gunshot or drawing a cork” whereas `pop.n.04` corresponds to “music of general appeal to teenagers; ...”. For a given pair of words  $(u, v)$ , we use the score of the synset pair  $(s'_u, s'_v)$  that has the lowest KL divergence among all the pairs  $S_u \times S_v$  where  $S_u, S_v$  are sets of synsets for words  $u$  and  $v$ , respectively. More precisely,  $s(u, v) = -\min_{s_u \in S_u, s_v \in S_v} D(s_u, s_v)$ . This pair selection corresponds to choosing the synset pair that has the highest degree of entailment. This approach has been used in word embeddings literature to select most related word pairs (Athiwaratkun & Wilson, 2017). For word pairs that are not in the model, we assign the score equal to the median of all scores. We evaluate our model scores against the human annotated scores using Spearman’s rank correlation.

Table 3 shows HYPERLEX results of our models **DOE-A** (asymmetric) and **DOE-S** (symmetric) as well as other competing models. The model **DOE-A** which uses KL divergence and negative sampling approach **S1**, **S2** and **S4** outperforms all other existing models, achieving state-of-the-art performance for the HYPERLEX noun dataset. (See Section A.1 for hyperparameter details) The model **DOE-S** which uses expected likelihood kernel attains a lower score of 0.455 compared to the asymmetric counterpart (**DOE-A**). This result underscores the importance of asymmetric measures which can capture relationship directionality.

We provide a brief summary of competing models: **FR** scores are based on concept word frequency ratio (Weeds et al., 2004). **SLQS** uses entropy-based measure to quantify entailment (Santus et al., 2014). **Vis-ID** calculates scores based on visual generality measures (Kiela et al., 2015). **WN-B** calculates the scores based on the shortest path between concepts in WN taxonomy (Miller, 1995). **w2g** Gaussian embeddings trained using the methodology in Vilnis & McCallum (2015). **VOE** Vector order embeddings (Vendrov et al., 2016). **Euc** and **Poin** calculate scores based on the Euclidean distance and Poincaré distance of the trained Poincaré embeddings (Nickel & Kiela, 2017). The models **FR** and **SLQS** are based on word occurrences in text corpus, where **FR** is trained on the British National Corpus and **SLQS** is trained on UKWAC, WACKYPEDIA (Bailey & Thompson, 2006; Baroni et al., 2009) and annotated BLESS dataset (Baroni & Lenci, 2011). Other models **Vis-ID**, **w2g**, **VOE**, **Euc**, **Poin** and ours are trained on WordNet, with the exception that **Vis-ID** also uses



Table 3: Spearman’s correlation for HYPERLEX nouns.

	FR	SLQS	Vis-ID	WN-B	w2g	VOE	Poin	HypV	DOE-S	DOE-A
$\rho$	0.283	0.229	0.253	0.240	0.192	0.195	0.512	0.540	0.455	<b>0.590</b>

Table 4: Spearman’s correlation for HYPERLEX nouns for different negative sample schemes.

Negative Samples	$\rho$	Negative Samples	$\rho$
$1 \times \mathbf{S1}$	0.527	$1 \times \mathbf{S1} + \mathbf{S2} + \mathbf{S4}$	<b>0.590</b>
$2 \times \mathbf{S1}$	0.529	$2 \times \mathbf{S1} + \mathbf{S2} + \mathbf{S4}$	0.580
$5 \times \mathbf{S1}$	0.518	$5 \times \mathbf{S1} + \mathbf{S2} + \mathbf{S4}$	0.582
$10 \times \mathbf{S1}$	0.517	$1 \times \mathbf{S1} + \mathbf{S2} + \mathbf{S3}$	0.570
$1 \times \mathbf{S1} + \mathbf{S2}$	0.567	$2 \times \mathbf{S1} + \mathbf{S2} + \mathbf{S3}$	0.581
$2 \times \mathbf{S1} + \mathbf{S2}$	0.567	$\mathbf{S1} + 0.1 \times \mathbf{S2} + 0.9 \times \mathbf{S3}$	0.564
$3 \times \mathbf{S1} + \mathbf{S2}$	0.584	$\mathbf{S1} + 0.3 \times \mathbf{S2} + 0.7 \times \mathbf{S3}$	0.574
$5 \times \mathbf{S1} + \mathbf{S2}$	0.561	$\mathbf{S1} + 0.7 \times \mathbf{S2} + 0.3 \times \mathbf{S3}$	0.555
$10 \times \mathbf{S1} + \mathbf{S2}$	0.550	$\mathbf{S1} + 0.9 \times \mathbf{S2} + 0.1 \times \mathbf{S3}$	0.533

Google image search results for visual data. The reported results of **FR**, **SLQS**, **Vis-ID**, **WN-B**, **w2g** and **VOE** are from Vulić et al. (2016).

We note that an implementation of Gaussian embeddings model (**w2g**) reported by Vulić et al. (2016) does not perform well compared to previous benchmarks such as **Vis-ID**, **FR**, **SLQS**. Our training approach yields the opposite results and outperforms other highly competitive methods such as Poincaré embeddings and Hypervec. This result highlights the importance of the training approach, even if the concept representation of our work and Vilnis & McCallum (2015) both use Gaussian distributions. In addition, we observe that vector order embeddings (VOE) do not perform well compared to our model, which we hypothesize is due to the “soft” orders induced by the divergence penalty that allows our model scores to more closely reflect hypernymy degrees.

We note another interesting observation that a model trained on a symmetric divergence (ELK) from Section 4.2 can also achieve a high HYPERLEX correlation of 0.532 if KL is used to calculate the model scores. This is because the encapsulation behavior can arise even though the training penalty is symmetric (more explanation in Section 4.2). However, using the symmetric divergence based on ELK results in poor performance on HYPERLEX (0.455), which is expected since it cannot capture the directionality of hypernymy.

We note that another model LEAR obtains an impressive score of 0.686 (Vulić & Mrkšić, 2014). However, LEAR use pre-trained word embeddings such as WORD2VEC or GLOVE as a pre-processing step, leveraging a large vocabulary with rich semantic information. To the best of our knowledge, our model achieves the highest HYPERLEX Spearman’s correlation among models without using large-scale pre-trained embeddings.

Table 4 shows the effects of negative sample selection described in Section 3.5. We note again that **S1** is the technique used in literature Socher et al. (2013); Vendrov et al. (2016) and **S2**, **S3**, **S4** are the new techniques we proposed. The notation, for instance,  $k \times \mathbf{S1} + \mathbf{S2}$  corresponds to using  $k$  samples from **S1** and 1 sample from **S2** per each positive sample. We observe that our new selection methods offer strong improvement from the range of 0.51 – 0.52 (using **S1** alone) to 0.55 or above for most combinations with our new selection schemes.

## 5 FUTURE WORK

Analogous to recent work by Vulić & Mrkšić (2014) which post-processed word embeddings such as GLOVE or WORD2VEC, our future work includes using the WordNet hierarchy to impose encapsulation orders when training probabilistic embeddings.

In the future, the distribution approach could also be developed for encoder-decoder based models for tasks such as caption generation where the encoder represents the data as a distribution, containing semantic and visual features with uncertainty, and passes this distribution to the decoder which maps to text or images. Such approaches would be reminiscent of variational autoencoders (Kingma & Welling, 2013), which take *samples* from the encoder’s distribution.

#### ACKNOWLEDGEMENTS

We thank NSF IIS-1563887 for support.

#### REFERENCES

- Ben Athiwaratkun and Andrew Gordon Wilson. Multimodal word distributions. In *ACL*, 2017.
- Steve Bailey and Dave Thompson. UKWAC: building the uk’s first public web archive. *D-Lib Magazine*, 12(1), 2006.
- Marco Baroni and Alessandro Lenci. How we blessed distributional semantic evaluation. In *Proceedings of the GEMS 2011 Workshop on GEometrical Models of Natural Language Semantics*, 2011.
- Marco Baroni, Silvia Bernardini, Adriano Ferraresi, and Eros Zanchetta. The wacky wide web: a collection of very large linguistically processed web-crawled corpora. *Language Resources and Evaluation*, 43(3), 2009.
- Marco Baroni, Raffaella Bernardi, Ngoc-Quynh Do, and Chung-chieh Shan. Entailment above the word level in distributional semantics. In *EACL*, pp. 23–32, 2012.
- Yonatan Bisk, Deniz Yuret, and Daniel Marcu. Natural language communication with robots. In *NAACL*, 2016.
- Julia Hockenmaier and Alice Lai. Learning to predict denotational probabilities for modeling entailment. In *EACL*, 2017.
- Tony Jebara, Risi Kondor, and Andrew Howard. Probability product kernels. *JMLR*, 2004.
- Douwe Kiela, Laura Rimell, Ivan Vulić, and Stephen Clark. Exploiting image generality for lexical entailment detection. In *ACL*, 2015.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *CoRR*, abs/1412.6980, 2014.
- Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *CoRR*, abs/1312.6114, 2013.
- Ryan Kiros, Yukun Zhu, Ruslan Salakhutdinov, Richard S. Zemel, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Skip-thought vectors. In *NIPS*, 2015.
- Xiang Li, Luke Vilnis, and Andrew McCallum. Improved representation learning for predicting commonsense ontologies. *CoRR*, abs/1708.00549, 2017. URL <http://arxiv.org/abs/1708.00549>.
- Yingzhen Li and Richard E. Turner. Rényi divergence variational inference. In *NIPS*, 2016.
- Friedrich Liese and Igor Vajda. *Convex Statistical Distances*. Leipzig : Teubner, 1987.
- Edward Loper and Steven Bird. Nltk: The natural language toolkit. In *ACL workshop*, 2002.
- Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *CoRR*, abs/1301.3781, 2013.
- George A. Miller. Wordnet: A lexical database for english. *Commun. ACM*, 38(11), November 1995.
- Maximilian Nickel and Douwe Kiela. Poincaré embeddings for learning hierarchical representations. *NIPS*, 2017.
- Leandro Pardo. *Statistical Inference Based on Divergence Measures*, chapter 1, pp. 1–54. Chapman & Hall/CRC, 2006.
- Alfred Renyi. On measures of entropy and information. In *Proceedings of the Fourth Berkeley Symposium on Mathematical Statistics and Probability, Volume 1: Contributions to the Theory of Statistics*, Berkeley, Calif., 1961.

- Tim Rocktäschel, Edward Grefenstette, Karl Moritz Hermann, Tomáš Kociský, and Phil Blunsom. Reasoning about entailment with neural attention. *CoRR*, abs/1509.06664, 2015.
- Enrico Santus, Alessandro Lenci, Qin Lu, and Sabine Schulte im Walde. Chasing hypernyms in vector spaces with entropy. In *EACL*, 2014.
- Richard Socher, Danqi Chen, Christopher D. Manning, and Andrew Y. Ng. Reasoning with neural tensor networks for knowledge base completion. In *NIPS*, 2013.
- Robert Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of general knowledge. *CoRR*, abs/1612.03975, 2016. URL <http://arxiv.org/abs/1612.03975>.
- Ivan Vendrov, Ryan Kiros, Sanja Fidler, and Raquel Urtasun. Order-embeddings of images and language. *ICLR*, 2016.
- Luke Vilnis and Andrew McCallum. Word representations via gaussian embedding. *ICLR*, 2015.
- Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015.
- Ivan Vulić and Nikola Mrkšić. Specialising word vectors for lexical entailment. *CoRR*, abs/1710.06371, 2014.
- Ivan Vulić, Daniela Gerz, Douwe Kiela, Felix Hill, and Anna Korhonen. Hyperlex: A large-scale evaluation of graded lexical entailment. *CoRR*, 2016.
- Julie Weeds, David J. Weir, and Diana McCarthy. Characterising measures of lexical distributional similarity. In *COLING*, 2004.

## A SUPPLEMENTARY MATERIALS

### A.1 MODEL HYPERPARAMETERS

In Section 4.3, the 2-dimensional Gaussian model is trained with **S-1** method where the number of negative samples is equal to the number of positive samples. The best hyperparameters for  $d = 2$  model is  $(m, \beta, \gamma) = (100.0, 2 \times 10^{-4}, 3.0)$ .

In Section 4.2, the best hyperparameters  $(m, \beta, \gamma)$  for each of our model are as follows: For Gaussian with KL penalty:  $(2000.0, 5 \times 10^{-5}, 500.0)$ , Gaussian with reversed KL penalty:  $(1000.0, 1 \times 10^{-4}, 1000.0)$ , Gaussian with ELK penalty  $(1000, 1 \times 10^{-5}, 10)$ .

In Section 4.4, we use the same hyperparameters as in 4.2 with KL penalty, but a different negative sample combination in order to increase the distinguishability of divergence scores. For each positive sample in the training set, we use one sample from each of the methods **S1**, **S2**, **S4**. We note that the model from Section 4.2, using **S1** with the KL penalty obtains a Spearman’s correlation of 0.527.

### A.2 ANALYSIS OF TRAINING METHODOLOGY

We emphasize that Gaussian embeddings have been used in the literature, both in the unsupervised settings where word embeddings are trained with local contexts from text corpus, and in supervised settings where concept embeddings are trained to model annotated data such as WORDNET . The results in supervised settings such as modeling WORDNET have been reported to compare with competing models but often have inferior performance (Vendrov et al., 2016; Vulić et al., 2016). Our paper reaches the opposite conclusion, showing that a different training approach using Gaussian representations can achieve state-of-the-art results.

#### A.2.1 DIVERGENCE THRESHOLD

Consider a relationship  $f \preceq g$  where  $f$  is a hyponym of  $g$  or  $g$  is a hypernym of  $f$ . Even though the divergence  $D(f||g)$  can capture the extent of encapsulation, a density  $f$  will have the lowest divergence with respect with  $g$  only if  $f = g$ . In addition, if  $f$  is a more concentrated distribution that is encompassed by  $g$ ,  $D(f||g)$  is minimized when  $f$  is at the center of  $g$ . However, if there are many hyponyms  $f_1, f_2$  of  $g$ , the hyponyms can compete to be close to the center, resulting in too much overlapping between  $f_1$  and  $f_2$  if the random sampling to penalize negative pairs is not sufficiently strong. The divergence threshold  $\gamma$  is used such that there is no longer a penalty once the divergence is below a certain level.

We demonstrate empirically that the threshold  $\gamma$  is important for learning meaningful Gaussian distributions. We fix the hyperparameters  $m = 2000$  and  $\beta = 5 \times 10^{-5}$ , with **S1** negative sampling. Figure 5 shows that there is an optimal non-zero threshold and yields the best performance for both WORDNET Hypernym prediction and HYPERLEX Spearman’s correlation. We observe that using  $\gamma = 0$  is detrimental to the performance, especially on HYPERLEX results.

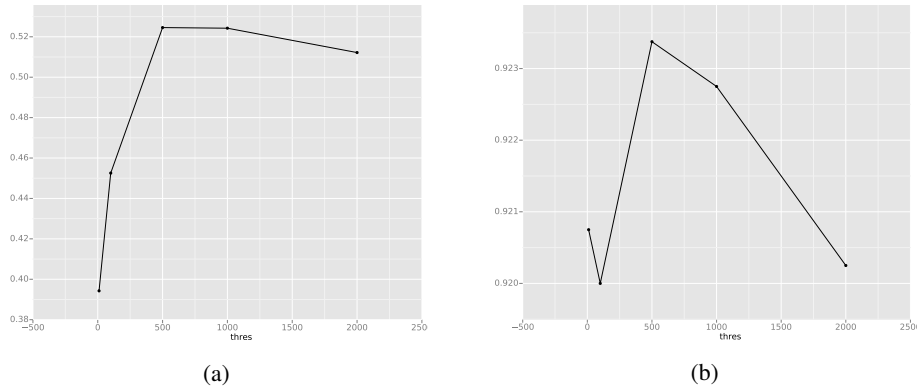


Figure 5: (a) Spearman’s correlation on HYPERLEX versus  $\gamma$  (b) Test Prediction Accuracy versus  $\gamma$ .

#### A.2.2 INITIAL VARIANCE SCALE

As opposed to the mean vectors that are randomly initialized, we initialize all diagonal covariance elements to be the same. Even though the variance can adapt during training, we find that different initial scales of variance result in drastically different performance. To demonstrate, in Figure 6, we show the best test accuracy and

HYPERLEX Spearman’s correlation for each initial variance scale, with other hyperparameters (margin  $m$  and threshold  $\gamma$ ) tuned for each variance. We use **S1 + S2 + S4** as a negative sampling method. In general, a low variance scale  $\beta$  increases the scale of the loss and requires higher margin  $m$  and threshold  $\gamma$ . We observe that the best prediction accuracy is obtained when  $\log(\beta) \approx -10$  or  $\beta = 5 \times 10^{-5}$ . The best HYPERLEX results are obtained when the scales of  $\beta$  are sufficiently low. The intuition is that low  $\beta$  increases the scale of divergence  $D(\cdot||\cdot)$ , which increases the ability to capture relationship nuances.

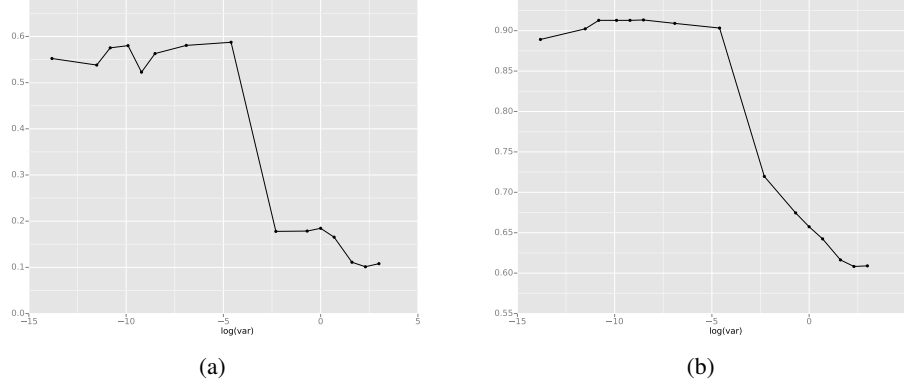


Figure 6: (a) Spearman’s correlation on HYPERLEX versus  $\log(\beta)$  (b) Test Prediction Accuracy versus  $\log(\beta)$ .

#### A.2.3 LOSS FUNCTION

We verify that for this task, our loss function in Equation 7 is superior to Equation 1 originally proposed by Vilnis & McCallum (2015). We use the exact same setup with new negative sample selections and KL divergence thresholding and compare the two loss functions. Table 5 verifies our claim.

Table 5: Best results for each loss function for two negative sampling setups: **S1 (Left)** and **S1 + S2 + S4 (Right)**

	Test Accuracy	HYPERLEX		Test Accuracy	HYPERLEX
Eq. 7	0.923	0.527	Eq. 7	0.911	0.590
Eq. 1	0.886	0.524	Eq. 1	0.796	0.489

#### A.2.4 DIMENSIONALITY

Table 6 shows the results for many dimensionalities for two negative sample strategies: **S1** and **S1 + S2 + S4**.

Table 6: Best results for each dimension with negative samples **S1 (Left)** and **S1 + S2 + S4 (Right)**

$d$	Test Accuracy	HYPERLEX	$d$	Test Accuracy	HYPERLEX
5	0.909	0.437	5	0.901	0.483
10	0.919	0.462	10	0.909	0.526
20	0.922	0.487	20	0.914	0.545
50	0.923	0.527	50	0.911	0.590
100	0.924	0.526	100	0.913	0.573
200	0.918	0.526	200	0.910	0.568

#### A.2.5 $\alpha$ -DIVERGENCES

Table 7 show the results using models trained and evaluated with  $D(\cdot||\cdot) = D_\alpha(\cdot||\cdot)$  with negative sampling approach **S1**. Interestingly, we found that  $\alpha \rightarrow 1$  (KL) offers the best result for both prediction accuracy and

**HYPERLEX**. It is possible that  $\alpha = 1$  is sufficiently asymmetric enough to distinguish hypernym directionality, but does not have as sharp penalty as in  $\alpha > 1$ , which can help learning.

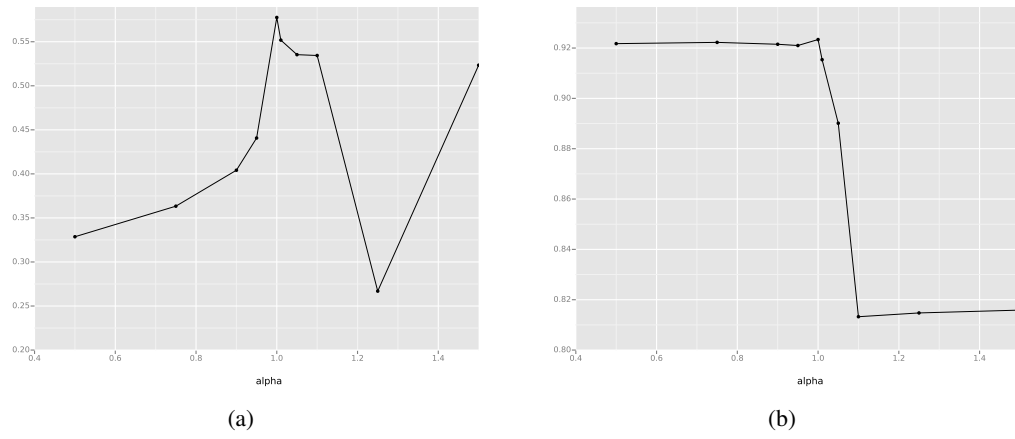


Figure 7: (a) Spearman's correlation on HYPERLEX versus  $\alpha$  (b) Test Prediction Accuracy versus  $\alpha$ .

# LOWER AND UPPER BOUNDS FOR APPROXIMATION OF THE KULLBACK-LEIBLER DIVERGENCE BETWEEN GAUSSIAN MIXTURE MODELS

*J.-L. Durrieu, J.-Ph. Thiran*

Signal Processing Laboratory (LTS5)  
École Polytechnique Fédérale de Lausanne (EPFL)  
Switzerland

*F. Kelly*

Dept. of Electronic and Electrical Engineering  
Trinity College Dublin  
Ireland

## ABSTRACT

Many speech technology systems rely on Gaussian Mixture Models (GMMs). The need for a comparison between two GMMs arises in applications such as speaker verification, model selection or parameter estimation. For this purpose, the Kullback-Leibler (KL) divergence is often used. However, since there is no closed form expression to compute it, it can only be approximated. We propose lower and upper bounds for the KL divergence, which lead to a new approximation and interesting insights into previously proposed approximations. An application to the comparison of speaker models also shows how such approximations can be used to validate assumptions on the models.

**Index Terms**— Gaussian Mixture Model (GMM), Kullback-Leibler Divergence, speaker comparison, speech processing.

## 1. INTRODUCTION

Gaussian Mixture Models (GMMs) are widely used to model unknown probability density functions (PDFs). GMMs have many properties that make them particularly useful for parameter estimation. Kullback-Leibler divergences between two PDFs  $f$  and  $g$ ,  $D_{\text{KL}}(f||g)$  can be used to compare such distributions. They arise in various (speech processing) applications: to classify speakers [1], as a cost to minimize for parameter estimation [2] or as a Kernel for Support Vector Machines (SVMs) [3, 4].

Let  $f$  and  $g$  be two PDFs, defined on  $\mathbb{R}^d$ , where  $d$  is the dimension of the observed vectors  $\mathbf{x}$ . The Kullback-Leibler divergence (KL divergence) between  $f$  and  $g$  is defined as:

$$D_{\text{KL}}(f||g) = \int_{\mathbb{R}^d} f(\mathbf{x}) \log \frac{f(\mathbf{x})}{g(\mathbf{x})} d\mathbf{x} \quad (1)$$

When  $f$  and  $g$  are the PDFs of normal random multivariate variables, *i.e.*

$$\log f(\mathbf{x}) = -\frac{1}{2} \log \left( (2\pi)^d |\Sigma^f| \right) - \frac{1}{2} (\mathbf{x} - \boldsymbol{\mu}^f)^T (\Sigma^f)^{-1} (\mathbf{x} - \boldsymbol{\mu}^f) \\ f(\mathbf{x}) \triangleq N(\mathbf{x}; \boldsymbol{\mu}^f, \Sigma^f) \text{ and } g(\mathbf{x}) \triangleq N(\mathbf{x}; \boldsymbol{\mu}^g, \Sigma^g) \quad (2)$$

where  $\boldsymbol{\mu}^f$  and  $\Sigma^f$  ( $\boldsymbol{\mu}^g$  and  $\Sigma^g$ , respectively) are the mean and covariance matrix of  $f$  (resp.  $g$ ),  $T$  is the transpose operator and  $|\Sigma^f|$

the determinant of  $\Sigma^f$ , then the KL divergence between  $f$  and  $g$  has a closed form expression [5]:

$$D_{\text{KL}}(f||g) = \frac{1}{2} \log \frac{|\Sigma^g|}{|\Sigma^f|} + \frac{1}{2} \text{Tr}((\Sigma^g)^{-1} \Sigma^f) \\ + \frac{1}{2} (\boldsymbol{\mu}^f - \boldsymbol{\mu}^g)^T (\Sigma^g)^{-1} (\boldsymbol{\mu}^f - \boldsymbol{\mu}^g) - \frac{d}{2} \quad (3)$$

For GMMs, however, the KL divergence does not have such a closed form expression. Letting  $f$  and  $g$  now be the PDFs for two GMMs, the expression of  $f$  becomes (with an analogous expression for  $g$ ):

$$f(\mathbf{x}) = \sum_{a=1}^A \omega_a^f f_a(\mathbf{x}) = \sum_{a=1}^A \omega_a^f N(\mathbf{x}; \boldsymbol{\mu}_a^f, \Sigma_a^f) \quad (4)$$

where  $A$  and  $B$  are the number of components of the GMM for  $f$  and  $g$ , respectively, and where  $f_a$  and  $g_b$ ,  $\forall a, b$ , are individual normal PDFs. It is possible to obtain an accurate approximation to the KL divergence between  $f$  and  $g$ , via Monte-Carlo estimations, but only at a great computational cost. Fast and reliable approximations for the KL divergence are therefore sought after [6, 7]. We propose the calculation of a lower and an upper bound for the KL divergence between two GMMs. The mean of these bounds then provides an approximation comparable to the approximations proposed by Hershey and Olsen [6]. These bounds are essential when one needs to minimize or maximize the KL divergence, since minimizing the upper bounds implies minimizing the divergence.

We first describe previous proposals for approximations of the KL divergence. Then the proposed lower and upper bounds are derived, with discussions about their interpretations. Finally, some numerical results and an application to speaker model comparison are presented.

## 2. APPROXIMATIONS TO THE KULLBACK-LEIBLER DIVERGENCE

In this section, we recall the approximations presented in [6].

### 2.1. Monte Carlo Estimation

The KL divergence can be approximated via Monte-Carlo (MC) estimation. It can indeed be expressed as the expectation of the logarithm of the ratio of  $f$  over  $g$ , under the PDF  $f$ . Let  $X$  be a (multivariate) random variable, with PDF  $f$ . Then, by definition:

$$D_{\text{KL}}(f||g) = E_X [\log (f(X)/g(X))] \quad (5)$$

This work was partly funded by the Swiss CTI agency, project n. 11359.1 PFES-ES, in collaboration with SpeedLingua SA, Lausanne, Switzerland, and partly funded by the Irish Research Council for Science, Engineering and Technology.



The MC methodology can therefore be applied to estimate such expectations, by the following steps:

1. Draw  $n$  independent samples  $\mathbf{x}_i$  from the PDF  $f$ ,
2. Compute  $D_{\text{MC},n}(f||g) = \frac{1}{n} \sum_i \log(f(\mathbf{x}_i)/g(\mathbf{x}_i))$ .

By the law of large numbers,  $D_{\text{MC},n}(f||g)$  converges to  $D_{\text{KL}}(f||g)$  as  $n$  tends to infinity. In this work, we chose to consider this MC approximation with  $n = 10^6$  as a reference.

## 2.2. Product of Gaussians Approximation

Hershey and Olsen proposed a decomposition which serves as basis for several of the approximations [6], including the ones proposed here. Let  $L_f(g) = E_X[\log g(X)]$ , where  $X \sim f$ . The KL divergence can then be decomposed as:

$$D_{\text{KL}}(f||g) = L_f(f) - L_f(g) \quad (6)$$

The “product of Gaussians” approximation,  $D_{\text{prod}}$ , is derived thanks to (6) and Jensen’s inequality to find upper bounds for  $L_f(g)$  and  $L_f(f)$ :

$$L_f(g) = \sum_a \omega_a^f \int_{\mathbf{x}} f_a(\mathbf{x}) \log\left(\sum_b \omega_b^g g_b(\mathbf{x})\right) d\mathbf{x} \quad (7)$$

$$\leq \sum_a \omega_a^f \log\left(\sum_b \omega_b^g \int_{\mathbf{x}} f_a(\mathbf{x}) g_b(\mathbf{x}) d\mathbf{x}\right) \quad (8)$$

$$L_f(g) \leq \sum_a \omega_a^f \log\left(\sum_b \omega_b^g t_{ab}\right) \quad (9)$$

where  $t_{ab} \triangleq \int_{\mathbf{x}} f_a(\mathbf{x}) g_b(\mathbf{x}) d\mathbf{x}$  is the normalization constant of the product of the Gaussians. Similarly, we have:

$$L_f(f) \leq \sum_a \omega_a^f \log\left(\sum_{\alpha} \omega_{\alpha}^f z_{a\alpha}\right) \quad (10)$$

$$z_{a\alpha} \triangleq \int_{\mathbf{x}} f_a(\mathbf{x}) f_{\alpha}(\mathbf{x}) d\mathbf{x} \quad (11)$$

Assuming that these upper bounds are close enough to  $L_f(g)$  and  $L_f(f)$ , respectively, these latter quantities can be approximated by their upper bounds, in order to derive  $D_{\text{prod}}$  [6]:

$$D_{\text{prod}}(f||g) \triangleq \sum_a \omega_a^f \log \frac{\sum_{\alpha} \omega_{\alpha}^f z_{a\alpha}}{\sum_b \omega_b^g t_{ab}} \quad (12)$$

The closed form expression of the normalization constants is given in Appendix A.

## 2.3. Variational Approximation

Lower bounds for  $L_f(g)$  and  $L_f(f)$  can also be derived, using variational parameters as follows [6]:

$$L_f(g) = E_X[\log(\sum_b \omega_b^g g_b(\mathbf{x}))] \quad (13)$$

$$= \sum_a \omega_a^f \int_{\mathbf{x}} f_a(\mathbf{x}) \log\left(\sum_b \omega_b^g \phi_{ba} \frac{g_b(\mathbf{x})}{\phi_{ba}}\right) d\mathbf{x} \quad (14)$$

$$\geq \sum_{ab} \omega_a^f \phi_{ba} \int_{\mathbf{x}} f_a(\mathbf{x}) \log \frac{\omega_b^g g_b(\mathbf{x})}{\phi_{ba}} d\mathbf{x} \quad (15)$$

where  $\phi_{ba} \geq 0$ , with  $\sum_b \phi_{ba} = 1, \forall a, b$ . Maximizing the right hand side of the above equation, with respect to  $\phi_{ba}$ , provides a lower bound to  $L_f(g)$ :

$$L_f(g) \geq \sum_a \omega_a^f \log \sum_b \omega_b^g e^{-D_{\text{KL}}(f_a||g_b)} - \sum_a \omega_a^f H(f_a) \quad (16)$$

where  $H(f_a)$  is the entropy of  $f_a$ , with a closed form given in Appendix B, and where  $D_{\text{KL}}(f_a||g_b)$  also has a closed form expression, as given in Eq. (3). Similarly,  $L_f(f)$  has the following variational lower bound:

$$L_f(f) \geq \sum_{\alpha} \omega_{\alpha}^f \log \sum_a \omega_a^f e^{-D_{\text{KL}}(f_{\alpha}||f_a)} - \sum_{\alpha} \omega_{\alpha}^f H(f_{\alpha}) \quad (17)$$

As in the previous section, these lower bounds can be used as approximations for the corresponding quantities in order to derive the “variational” approximation [6]:

$$D_{\text{var}}(f||g) = \sum_a \omega_a^f \log \frac{\sum_{\alpha} \omega_{\alpha}^f e^{-D_{\text{KL}}(f_{\alpha}||f_a)}}{\sum_b \omega_b^g e^{-D_{\text{KL}}(f_a||g_b)}} \quad (18)$$

These simple closed form expressions make it easy to compute an approximation to  $D_{\text{KL}}$ , with properties close to that of  $D_{\text{KL}}$ . However, there does not seem to be a theoretical reason why these quantities should be approximations to  $D_{\text{KL}}$ , although numerical results have shown their relevance [6]. Since  $D_{\text{prod}}$  and  $D_{\text{var}}$  are each the sum of an upper bound with a lower bound, it is difficult to analyze in what sense they approximate the KL divergence.

Based on similar principles, we propose upper and lower bounds that shed a new light on these approximations.

## 3. UPPER AND LOWER BOUNDS FOR THE KL DIVERGENCE

Strict bounds are mainly useful in the parameter estimation case, and by providing the interval in which we can find the real value of the KL divergence, they provide a well motivated way to design another approximation to the divergence. Using the KL decomposition (6) and the above individual bounds, we propose the following bounds:

*Lower bound:* Combining Eqs. (9) and (17), we obtain the following lower bound for the KL divergence between GMMs:

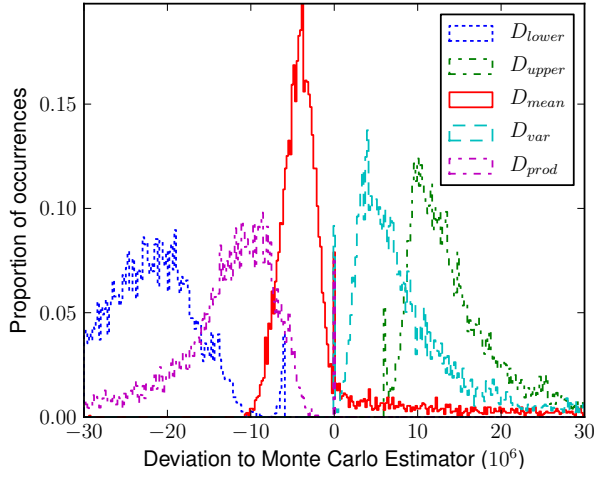
$$\underbrace{\sum_a \omega_a^f \log \frac{\sum_{\alpha} \omega_{\alpha}^f e^{-D_{\text{KL}}(f_{\alpha}||f_a)}}{\sum_b \omega_b^g t_{ab}} - \sum_a \omega_a^f H(f_a)}_{D_{\text{lower}}(f||g)} \leq D_{\text{KL}}(f||g) \quad (19)$$

*Upper bound:* Similarly, from Eqs. (10) and (16), we obtain:

$$D_{\text{KL}}(f||g) \leq \underbrace{\sum_a \omega_a^f \log \frac{\sum_{\alpha} \omega_{\alpha}^f z_{a\alpha}}{\sum_b \omega_b^g e^{-D_{\text{KL}}(f_a||g_b)}} + \sum_a \omega_a^f H(f_a)}_{D_{\text{upper}}(f||g)} \quad (20)$$

It is worth calculating the mean of  $D_{\text{lower}}$  and  $D_{\text{upper}}$ , the “center” of the interval. This is in fact equal to the mean of  $D_{\text{prod}}$  and  $D_{\text{var}}$ :

$$\begin{aligned} D_{\text{mean}}(f||g) &\triangleq [D_{\text{upper}}(f||g) + D_{\text{lower}}(f||g)]/2 \\ &= [D_{\text{prod}}(f||g) + D_{\text{var}}(f||g)]/2 \end{aligned} \quad (21)$$



**Fig. 1.** Histograms of the approximation deviations to the MC estimator,  $d = 39$ .

Since this value is between the lower and upper bounds of the KL divergence, it is a KL approximation as reasonable as  $D_{\text{prod}}$  or  $D_{\text{var}}$ . Eq. (21) provides some insight into the results given in [6]: the authors noticed therein that  $D_{\text{prod}}$  tended to greatly underestimate  $D_{\text{KL}}$ , while  $D_{\text{var}}$  was among the best choices as an approximation for  $D_{\text{KL}}$ . The relation (21) helps us understand why these values can also be considered as approximations, even though their definitions in [6] do not allow much interpretation.

One should also note that for a Gaussian PDF  $f$ ,  $D_{\text{upper}}(f||f) = -D_{\text{lower}}(f||f) = \frac{d}{2}(1 - \log 2)$ . These “limits”, which appear also for GMMs, reveal that the proposed bounds may not be as tight as desired, in spite of the tighter “variational” part of the bound. However, their mean in this case is 0, and  $D_{\text{mean}}$  is therefore not influenced by these limits. Of the 3 properties of the KL divergence in [6],  $D_{\text{mean}}$ , like  $D_{\text{prod}}$  and  $D_{\text{var}}$ , satisfies the similarity property but not those of identifiability or positivity.

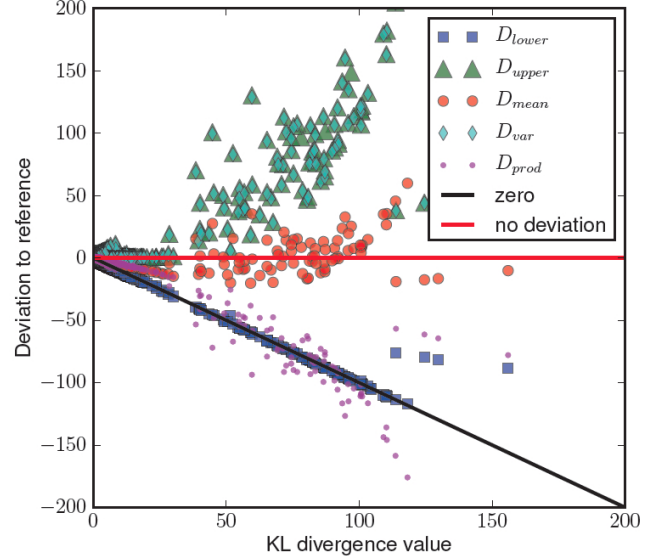
Finally, one should note that the complexities of the different approximations and bounds are roughly equivalent, in  $\mathcal{O}(K^2d)$  for diagonal covariance matrices and equal number of GMM components  $K$ . For the MC estimation, the complexity is in  $\mathcal{O}(NKn)$ . Since obtaining a reliable MC estimation requires  $N \gg K$ , the use of approximations is clearly advantageous from the computational complexity aspect.

## 4. NUMERICAL SIMULATIONS AND DISCUSSIONS

### 4.1. Deviation analysis

In order to compare these bounds and approximations, we created 100 synthetic GMMs, with the number of components  $K$  varying from 1 to 10 (10 GMMs for each value of  $K$ ), for each of the following dimensions  $d$  for the vectors: 1, 3, 39. The deviations of the approximations and bounds to the MC estimator of  $D_{\text{KL}}$ , with  $n = 10^6$  as the reference, are analyzed.

The histograms of the deviations for the different approximations and bounds are shown on Fig. 1, for  $d = 39$ . As expected,  $D_{\text{lower}}$  and  $D_{\text{upper}}$  are respectively below and above the reference. They however tend to greatly under- and over-estimate  $D_{\text{KL}}$ . They



**Fig. 2.** Deviations from the MC estimator against the reference KL divergence,  $d = 3$ . In addition to the quantities presented in the article, 2 lines represent the deviation of an “approximation” always equal to 0, and the “no deviation” line.

are therefore not suitable approximations to the desired divergence, specifically  $D_{\text{lower}}$  which is actually almost always close to 0, as can be seen on Fig. 2.

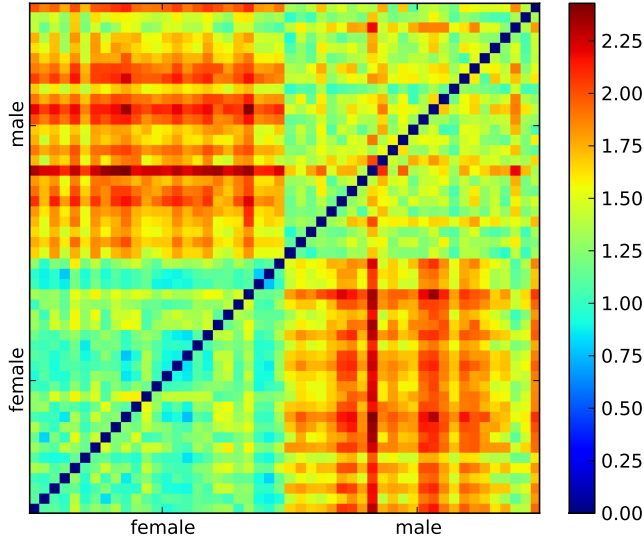
$D_{\text{var}}$  and  $D_{\text{prod}}$  are usually closer to  $D_{\text{KL}}$ , but, as expected, there is no rule as whether they are above or under  $D_{\text{KL}}$ : for  $d = 1$  and  $d = 3$ , the corresponding histograms even overlap.  $D_{\text{prod}}$  is generally under  $D_{\text{KL}}$ , while  $D_{\text{var}}$  slightly over-estimates it.  $D_{\text{mean}}$  seems to be closer to the desired value, with deviations more concentrated near 0. According to Fig. 2, the choice of an approximation may also depend on the actual value of the divergence; for small divergences, the approximations appear to be equivalent. For higher values,  $D_{\text{mean}}$  is a closer fit to the divergence than  $D_{\text{var}}$ , which tends to overestimate  $D_{\text{KL}}$ .

### 4.2. Speaker model comparison

As mentioned, approximations to the KL divergence and its bounds have numerous applications in speech processing. One application is that of speaker comparison, where it can be used as a similarity measure between GMMs representing speakers [1]. We have carried out a speaker comparison using the derived bounds to illustrate this application.

GMMs were trained for 50 speakers (25 male, 25 female) from the YOHO [8] database via adaptation of a gender-independent Universal Background Model (UBM) of 512 mixtures using 5 minutes of data [9]. Pre-processing involved energy-based silence removal and extraction of MFCC vectors of length 12 appended with delta and acceleration coefficients. The 50 models were compared by extracting  $D_{\text{mean}}$  between each model pair.

A confusion matrix of the comparisons is given in Fig. 3. The clusters of the within-gender and between-gender comparisons are easily identifiable. Between-gender divergence is generally greater than within-gender. This aligns with intuitive expectations about the relationship between male and female speaker models in the acoustic



**Fig. 3.** Confusion matrix for model comparisons with  $D_{\text{mean}}$ ,  $d = 36$ ,  $K = 512$ .

space *i.e.* that male models are closer to one another than to female models.

By observation, the KL divergence approximation provides a good estimation of the separation of the real, large GMMs in this test case. However, further work is needed to quantify and directly compare the quality of the estimations in the case of real data.

Finally it is worth noting that the correlation between  $D_{\text{mean}}$  and  $D_{\text{var}}$  is very high, meaning that either could be used for comparison purposes.

## 5. CONCLUSIONS

In this article, a lower and an upper bound for the Kullback-Leibler divergence between two GMM PDFs are proposed. The mean of these bounds provides an approximation to the KL divergence which is shown to be equivalent to previously proposed approximations, with a clearer theoretical motivation.

The closed form expressions of the bounds can be used for model comparisons, model validation, classification, or even to compute gradients whenever KL divergences are involved, for parameter estimation, for instance. Using a similar principle as proposed here, it could also be possible to speed up Monte-Carlo approximations, as shown in [10].

The proposed results could be easily extended to any mixture model, with arbitrary distribution PDFs, provided that closed form expressions for individual PDF divergence exist. The proposed bounds and approximation could at last be extended to the case of hidden Markov models.

## A. PRODUCT OF TWO GAUSSIANS

The normalizing constant for the product of two normal PDFs  $f_a$  and  $g_b$  is given by [11]:

$$\log t_{ab} = -\frac{d}{2} \log 2\pi - \frac{1}{2} \log |\Sigma_a^f + \Sigma_b^g| - \frac{1}{2} (\mu_b^g - \mu_a^f)^T (\Sigma_a^f + \Sigma_b^g)^{-1} (\mu_b^g - \mu_a^f) \quad (22)$$

## B. ENTROPY OF A MULTIVARIATE NORMAL DISTRIBUTION

Let  $f$  be a multivariate normal PDF,  $f(\mathbf{x}) = N(\mathbf{x}; \mu, \Sigma)$ , where  $\mathbf{x} \in \mathbb{R}^d$ . The entropy  $H(f)$  of  $f$  is:

$$H(f) \triangleq - \int_{\mathbf{x}} f(\mathbf{x}) \log f(\mathbf{x}) d\mathbf{x} = \frac{1}{2} \log \left( (2\pi e)^d |\Sigma| \right) \quad (23)$$

## C. REFERENCES

- [1] M. Ben, M. Betser, F. Bimbot, and G. Gravier, "Speaker diarization using bottom-up clustering based on a parameter-derived distance between adapted GMMs," in *Proc. of International Conference on Spoken Language Processing*, Jeju Island, Korea, October 4-8 2004.
- [2] Z. Ghahramani and M.I. Jordan, "Factorial hidden Markov models," *Machine Learning*, vol. 29, no. 2, pp. 245–273, 1997.
- [3] P.J. Moreno and P.P. Ho, "A new SVM approach to speaker identification and verification using probabilistic distance kernels," in *Proc. of European Conference on Speech Communication and Technology*, Geneva, Switzerland, September 1-4 2003, vol. 3, pp. 2965–2968.
- [4] W.M. Campbell, D.E. Sturim, and D.A. Reynolds, "Support vector machines using GMM supervectors for speaker verification," *IEEE Signal Processing Letters*, vol. 13, no. 5, pp. 308–311, 2006.
- [5] S.J. Roberts and W.D. Penny, "Variational Bayes for generalized autoregressive models," Tech. Rep., Oxford University, May 22 2002.
- [6] J. R. Hershey and P. A. Olsen, "Approximating the Kullback Leibler divergence between Gaussian Mixture Models," in *Proc. of the International Conference on Audio, Speech and Signal Processing*, Honolulu, Hawaii, USA, April 15-20 2007, vol. 4, pp. IV–317.
- [7] W. M. Campbell and Z. N. Karam, "Simple and efficient speaker comparison using approximate KL divergence," in *Proc. of Interspeech*, Makuhari, Chiba, Japan, Sept. 26-30 2010, pp. 362 – 365.
- [8] J. Campbell and A. Higgins, "YOHO speaker verification," Linguistic Data Consortium, 1994.
- [9] D.A. Reynolds, T.F. Quatieri, and R.B. Dunn, "Speaker verification using adapted Gaussian mixture models," *Digital signal processing*, vol. 10, no. 1-3, pp. 19–41, 2000.
- [10] J.-Y. Chen, J. R. Hershey, P. A. Olsen, and E. Yashchin, "Accelerated Monte Carlo for Kullback-Leibler divergence between Gaussian mixture models," in *Proc. of the International Conference on Acoustics, Speech and Signal Processing*, Las Vegas, Nevada, USA, March 31-April 4 2008.
- [11] P. Ahrendt, "The multivariate Gaussian probability distribution," Tech. Rep., Technical University of Denmark, Jan. 2005.

# SUMMARY

**Multimodal Word Distributions**

**ACL 2017**

## **Introduction:**

This paper introduces multimodal word embeddings where each word is represented as a mixture of gaussian mixtures and each gaussian refers to a different meaning of the word in case of polysemy words.

## **Previous Approaches:**

**Word2vec** - In this approach, each word is represented as a single vector where words with similar meanings are nearer to one another.

**Single Gaussian Representation** - In this approach, each word is represented as a gaussian distribution where mean and covariance are learned from the data. This provides much richer representation than point vectors but it cannot handle polysemies.

To overcome the problem in the above approaches, this paper proposes represent each word as a mixture of gaussians where each gaussian corresponds to different meanings in case of polysemies.

## **Approach:**

Each word  $w$  in a dictionary is represented as a gaussian mixture with  $k$  components and the distribution of the word is given by:

The main goal is to learn the model parameters from a corpus of natural sentences.

Each mean vector of the gaussian can represent one meaning of a word. They use a maximum margin energy-based ranking objective function which is given by:

$$L_{\theta}(w, c, c') = \max(0, m - \log E_{\theta}(w, c) + \log E_{\theta}(w, c'))$$

and the energy function is given by -

$$E(f, g) = \int f(x)g(x) dx = \langle f, g \rangle_{L_2}$$

$$\log E_{\theta}(f, g) = \log \sum_{j=1}^K \sum_{i=1}^K p_i q_j e^{\xi_{i,j}}$$

$$\begin{aligned} \xi_{i,j} &\equiv \log \mathcal{N}(0; \vec{\mu}_{f,i} - \vec{\mu}_{g,j}, \Sigma_{f,i} + \Sigma_{g,j}) \\ &= -\frac{1}{2} \log \det(\Sigma_{f,i} + \Sigma_{g,j}) - \frac{D}{2} \log(2\pi) \\ &\quad - \frac{1}{2} (\vec{\mu}_{f,i} - \vec{\mu}_{g,j})^{\top} (\Sigma_{f,i} + \Sigma_{g,j})^{-1} (\vec{\mu}_{f,i} - \vec{\mu}_{g,j}) \end{aligned}$$

Here  $c$  is nearby word to  $w$  within a context window of  $l$ , and  $c'$  is a negative context word.

It is based on the concept that the words nearer to one another are similar. This loss function tries to push the similarity of a word and its positive context higher than the similarity of a word and its negative context by a margin of  $m$ .

The energy function captures similarity between  $i$ th component of word  $w_f$  and  $j$ th component of word  $w_g$ . So higher the similarity between these values, higher the energy value and lower will be the loss function.

## Experiments:

They use a concatenation of 2 datasets, UKWAC (2.5 billion tokens) and Wackypedia (1 billion tokens). For each gaussian mixture component of a word the mean represents the embedding of the  $i$ th component and the variance represents its uncertainty.  $D=50$ ,  $K=2$  and  $l=10$  for the experiments.

Since each word has multiple components the similarity scores that are used between two words are:

1. Expected Likelihood Kernel - Inner product between gaussian mixtures
2. Maximum Cosine Similarity - The maximum of cosine similarity between all pairs of components of two words  $w_f$  and  $w_g$ .
3. Minimum Euclidean Distance - The minimum of euclidean distance between all pairs of components of two words  $w_f$  and  $w_g$ .

The word embeddings are evaluated on several word similarity datasets and spearman correlation is calculated between the labels and the scores. The model w2gm outperforms all the other state-of-art models using various similarity measures.

Using a mixture of gaussians succeeds in modeling word uncertainty better than unimodal approaches as it has reduced variances for all components rather than high variance for a single gaussian in case of polysemous words.

## Conclusion:

The multimodal word representation introduced in this paper is successful in capturing different semantics of polysemous words, uncertainty, and entailment, and also perform favorably on word similarity benchmarks.

# Probabilistic FastText for Multi-Sense Word Embeddings

## ACL 2018

### Introduction:

This paper introduces word embeddings with multiple gaussian mixtures where mean of each mixture component is given by sum of n-grams. This model can handle not only words with multiple meanings but also rare misspelt or even unseen words.

### Previous Approaches:

Most of the word embeddings are based on dictionary-level embeddings and cannot handle rare/unseen words. To overcome this problem, FASTTEXT was introduced which used character-level embeddings where each word is modeled as a sum of n-gram vectors.

### Approach:

Each word is represented as a gaussian mixture.

Each component of the mixture can represent different word senses, and the mean vectors of each component decompose into vectors of n-grams, to capture character-level information. This embedding has the ability to handle rare words and even foreign polysemies with an improvement of 1% over the w2gm model on SCWS.

The overall approach is similar to that of w2gm except to the mean vector calculation in the gaussian mixture. The mean vectors in w2gm are dictionary level which leads to poor semantic estimate for rare words. To overcome this problem, using subword structures we can have the mean as :

$$\mu_w = \frac{1}{|NG_w| + 1} \left( v_w + \sum_{g \in NG_w} z_g \right)$$



**Experiments:**

The probabilistic FASTTEXT which combines subword structure with probabilistic embeddings is trained on UKWAC and Wackypedia datasets for English and for foreign languages it is trained on French, German, and Italian text corpuses. Parameters are fixed at  $K=2$ ,  $l=10$  and  $n=3,4,5$ .

It is observed that the subword embeddings prefer words with overlapping characters as nearest neighbors. This model is seen to outperform all state of art models including w2gm by 3.1% and FASTTEXT by 1.2%. Even in case of foreign languages, this model outperforms the others.

**Conclusion:**

This paper thus proposes the first ever model to handle polysemies rare and unseen words.

**Future work lies in:**

1. Trade-off between learning full covariance matrices for each word distribution, computational complexity, and performance as in this model we are directly using spherical covariance matrices.
2. Co-training on many languages.