

# Modeling Annotated Data

David M. Blei  
Division of Computer Science  
University of California, Berkeley  
Berkeley, CA 94720

Michael I. Jordan  
Division of Computer Science  
and Department of Statistics  
University of California, Berkeley  
Berkeley, CA 94720

## ABSTRACT

We consider the problem of modeling annotated data—data with multiple types where the instance of one type (such as a caption) serves as a description of the other type (such as an image). We describe three hierarchical probabilistic mixture models which aim to describe such data, culminating in *correspondence latent Dirichlet allocation*, a latent variable model that is effective at modeling the joint distribution of both types and the conditional distribution of the annotation given the primary type. We conduct experiments on the Corel database of images and captions, assessing performance in terms of held-out likelihood, automatic annotation, and text-based image retrieval.

## Categories and Subject Descriptors

G.3 [Mathematics of Computing]: Probability and Statistics—*statistical computing, multivariate statistics*

## General Terms

algorithms, experimentation

## Keywords

probabilistic graphical models, empirical Bayes, variational methods, automatic image annotation, image retrieval

## 1. INTRODUCTION

Traditional methods of information retrieval are organized around the representation and processing of a document in a (high-dimensional) word-space. Modern multimedia documents, however, are not merely collections of words, but can be collections of related text, images, audio, and cross-references. When working with a corpus of such documents, there is much to be gained from representations which can explicitly model associations among the different types of data.

In this paper, we consider probabilistic models for documents that consist of pairs of data streams. Our focus is on

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

SIGIR '03, July 28–August 1, 2003, Toronto, Canada.  
Copyright 2003 ACM 1-58113-646-3/03/0007 ...\$5.00.

problems in which one data type can be viewed as an *annotation* of the other data type. Examples of such data include images and their captions, papers and their bibliographies, and genes and their functions. In addition to the traditional goals of retrieval, clustering, and classification, annotated data lends itself to tasks such as automatic data annotation and retrieval of unannotated data from annotation-type queries.

A number of recent papers have considered generative probabilistic models for such multi-type or *relational* data [2, 6, 4, 13]. These papers have generally focused on models that jointly cluster the different data types, basing the clustering on latent variable representations that capture low-dimensional probabilistic relationships among interacting sets of variables.

In many annotation problems, however, the overall goal appears to be that of finding a *conditional* relationship between types, and in such cases improved performance may be found in methods with a more discriminative flavor. In particular, the task of annotating an unannotated image can be viewed formally as a classification problem—for each word in the vocabulary we must make a yes/no decision. Standard discriminative classification methods, however, generally make little attempt to uncover the probabilistic structure of either the input domain or the output domain. This seems ill-advised in the image/word setting—surely there are relationships among the words labeling an image, and these relationships reflect corresponding relationships among the regions in that image. Moreover, it seems likely that capturing these relationships would be helpful in annotating new images. With these issues in mind, we approach the annotation problem within a framework that exploits the best of both the generative and the discriminative traditions.

In this paper, we build a set of increasingly sophisticated models for a database of annotated images, culminating in *correspondence latent Dirichlet allocation* (CORR-LDA), a model that finds conditional relationships between latent variable representations of sets of image regions and sets of words. We show that, in this class of models, only CORR-LDA succeeds in providing both an excellent fit of the joint data and an effective conditional model of the caption given an image. We demonstrate its use in automatic image annotation, automatic region annotation, and text-based image retrieval.

## 2. IMAGE/CAPTION DATA

Our work has focused on images and their captions from the Corel database. Following previous work [2], each im-

age is segmented into regions by the N-cuts algorithm [12]. For each region, we compute a set of real-valued features representing visual properties such as size, position, color, texture, and shape. Each image and its corresponding caption is represented as a pair  $(\mathbf{r}, \mathbf{w})$ . The first element  $\mathbf{r} = \{r_1, \dots, r_N\}$  is a collection of  $N$  feature vectors associated with the regions of the image. The second element  $\mathbf{w} = \{w_1, \dots, w_M\}$  is the collection of  $M$  words of the caption.

We consider hierarchical probabilistic models of image/caption data which involve mixtures over underlying discrete and continuous variables. Conditional on the values of these latent variables, the region feature vectors are assumed to be distributed as a multivariate Gaussian distribution with diagonal covariance, and the caption words are assumed to be distributed as a multinomial distribution over the vocabulary.

We are interested in models that can perform three tasks: modeling the joint distribution of an image and its caption, modeling the conditional distribution of words given an image, and modeling the conditional distribution of words given a particular region of an image. The first task can be useful for clustering and organizing a large database of images. The second task is useful for automatic image annotation and text-based image retrieval. The third task is useful for automatically labeling and identifying a particular region of an image.

### 3. GENERATIVE MODELS

In this section, we describe two hierarchical mixture models of image/caption data, noting their strengths and limitations with respect to the three tasks described above. We introduce a third model, *correspondence latent Dirichlet allocation*, in which the underlying probabilistic assumptions are appropriate for all three tasks.

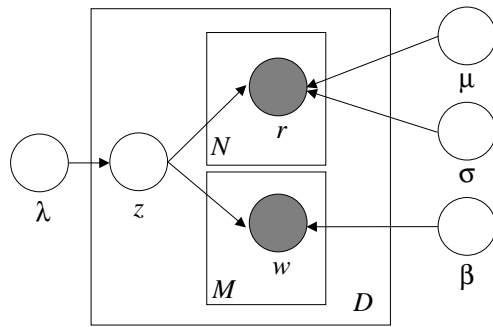
#### 3.1 A Gaussian-multinomial mixture model

We begin by considering a simple finite mixture model—the model underlying most previous work on the probabilistic modeling of multi-type data [2, 13]. In this model—the *Gaussian-multinomial mixture* (GM-MIXTURE) shown in Figure 1—a single discrete latent variable  $z$  is used to represent a joint clustering of an image and its caption. As shown in the figure, an image/caption is assumed to be generated by first choosing a value of  $z$ , and then repeatedly sampling  $N$  region descriptions  $r_n$  and  $M$  caption words  $w_m$  conditional on the chosen value of  $z$ . The variable  $z$  is sampled once per image/caption, and is held fixed during the process of generating its components. The joint distribution of the hidden factor  $z$  and the image/caption  $(\mathbf{r}, \mathbf{w})$  is:

$$p(z, \mathbf{r}, \mathbf{w}) = p(z | \lambda) \prod_{n=1}^N p(r_n | z, \mu, \sigma) \cdot \prod_{m=1}^M p(w_m | z, \beta). \quad (1)$$

Given a fixed number of factors  $K$  and a corpus of images/captions, the parameters of a GM-MIXTURE model can be estimated by the EM algorithm<sup>1</sup>. This yields  $K$  Gaussian distributions over features and  $K$  multinomial distributions over words which together describe a clustering of the images/captions. Since each image and its caption are

<sup>1</sup>In Section 4.2.1, we consider Bayesian versions of all of our models in which some the parameters are endowed with prior distributions.



**Figure 1: The GM-Mixture model of images and captions.** Following the standard graphical model formalism, nodes represent random variables and edges indicate possible dependence. Shaded nodes are observed random variables; unshaded nodes are latent random variables. The joint distribution can be obtained from the graph by taking the product of the conditional distribution of nodes given their parents (see Eq. 1). Finally, the box around a random variable is a “plate,” a notational device to denote replication. The box around  $r$  denotes  $N$  replicates of  $r$  (this gives the first product in Eq. 1).

assumed to have been generated conditional on the same factor, the resulting multinomial and Gaussian parameters will correspond. An image with high probability under a certain factor will likely contain a caption with high probability in the same factor.

Let us consider the three tasks under this model. First, the joint probability of an image/caption can be computed by simply marginalizing out the hidden factor  $z$  from Eq. (1). Second, we can obtain the conditional distribution of words given an image by invoking Bayes’ rule to find  $p(z | \mathbf{r})$  and marginalizing out the hidden factor:

$$p(w | \mathbf{r}) = \sum_z p(z | \mathbf{r}) p(w | z).$$

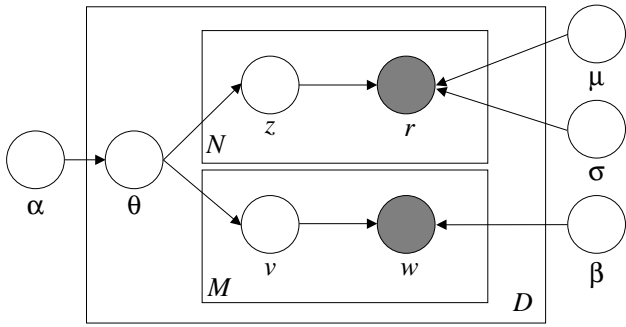
Finally, we would like to compute a region-specific distribution over words. This task, however, is beyond the scope of the GM-MIXTURE model. Conditional on the latent factor variable  $z$ , regions and words are generated independently, and the correspondence between specific regions and specific words is necessarily ignored.

#### 3.2 Gaussian-Multinomial LDA

The *latent Dirichlet allocation* (LDA) model is a latent variable model that allows factors to be allocated repeatedly within a given document or image [3]. Thus, different words in a document or different regions in an image can come from different underlying factors, and the document or image as a whole can be viewed as containing multiple “topics.” This idea can be applied directly to the joint modeling of images and captions.

*Gaussian-multinomial LDA* (GM-LDA), shown in Figure 2, assumes the following generative process:

1. Sample a Dirichlet random variable  $\theta$ ; this provides a probability distribution over the latent factors,  $\text{Mult}(\theta)$ .
2. For each of the  $N$  image regions:
  - (a) Sample  $z_n \sim \text{Mult}(\theta)$ .



**Figure 2: The GM-LDA model of images and captions. Unlike GM-Mixture (Figure 1), each word and image region can potentially be drawn from a different hidden factor.**

- (b) Sample a region description  $r_n$  conditional on  $z_n$ .
3. For each of the  $M$  words:
  - (a) Sample  $v_m \sim \text{Mult}(\theta)$ .
  - (b) Sample  $w_m$  conditional on  $v_m$ .

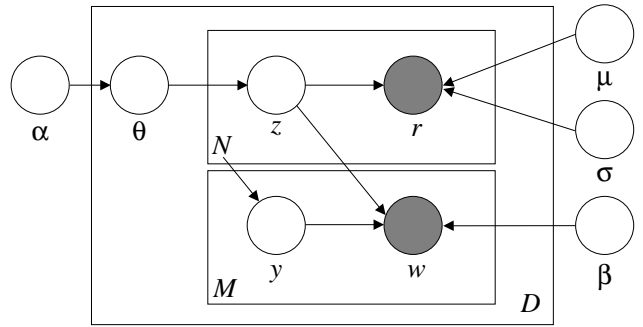
Note the importance of the plates. Within an image, all the region descriptions and words are generated with  $\theta$  held fixed; the latent factors for each word and region description can (potentially) vary. Each new image/caption is generated by again selecting from the Dirichlet variable  $\theta$  and repeating the entire process. Thus, we can view  $\theta$  as a high-level representation of the ensemble of image/caption pairs in terms of a probability distribution over factors that each image/caption can be assembled from.

The resulting joint distribution on image regions, caption words, and latent variables is given as follows:

$$p(\mathbf{r}, \mathbf{w}, \theta, \mathbf{z}, \mathbf{v}) = p(\theta | \alpha) \left( \prod_{n=1}^N p(z_n | \theta) p(r_n | z_n, \mu, \sigma) \right) \cdot \left( \prod_{m=1}^M p(v_m | \theta) p(w_m | v_m, \beta) \right).$$

As in the simpler LDA model, it is intractable to compute the conditional distributions of latent variables given observed data under this joint distribution, but efficient variational inference methods are available to compute approximations to these conditionals (see [2] for details). Furthermore, we can use variational inference methods to find the conditional probability  $p(w | \mathbf{r})$  needed for image annotation/retrieval and the conditional probability  $p(w | \mathbf{r}, r_n)$  needed for region labeling.

LDA provides significant improvements in predictive performance over simpler mixture models in the domain of text data [3], and we expect for GM-LDA to provide similar advantages over GM-MIXTURE. Indeed, we will see in Section 5 that GM-LDA does model the image/caption data better than GM-MIXTURE. We will also see, however, that good models of the joint probability of images and captions do not necessarily yield good models of the *conditional* probabilities that are needed for automatic annotation, text-based image retrieval, and region labeling. We will argue that this is due to the lack of a dependency between the latent factors  $z_n$  and  $v_m$  which respectively generated the images and their captions. In the next section, we turn to a model that aims to correct this problem.



**Figure 3: The graphical model representation of the Corr-LDA model. Note that the variables  $y_m$  are conditioned on  $N$ , the number of image regions.**

### 3.3 Correspondence LDA

We introduce *correspondence LDA* (CORR-LDA) as a model that combines the flexibility of GM-LDA with the associability of GM-MIXTURE. With this model, we achieve simultaneous dimensionality reduction in the representation of region descriptions and words, while also modeling the conditional correspondence between their respective reduced representations.

CORR-LDA is depicted in Figure 3. The model can be viewed in terms of a generative process that first generates the region descriptions and subsequently generates the caption words. In particular, we first generate  $N$  region descriptions  $r_n$  from an LDA model. Then, for each of the  $M$  caption words, one of the regions is selected from the image and a corresponding caption word  $w_m$  is drawn, conditioned on the factor that generated the selected region.

Formally, let  $\mathbf{z} = \{z_1, z_2, \dots, z_N\}$  be the latent factors that generate the image, and let  $\mathbf{y} = \{y_1, y_2, \dots, y_M\}$  be discrete indexing variables that take values from 1 to  $N$  with equal probability. Conditioned on  $N$  and  $M$ , a  $K$ -factor CORR-LDA model assumes the following generative process for an image/caption  $(\mathbf{r}, \mathbf{w})$ :

1. Sample  $\theta \sim \text{Dir}(\theta | \alpha)$ .
2. For each image region  $r_n$ ,  $n \in \{1, \dots, N\}$ :
  - (a) Sample  $z_n \sim \text{Mult}(\theta)$
  - (b) Sample  $r_n \sim p(r | z_n, \mu, \sigma)$  from a multivariate Gaussian distribution conditioned on  $z_n$ .
3. For each caption word  $w_m$ ,  $m \in \{1, \dots, M\}$ :
  - (a) Sample  $y_m \sim \text{Unif}(1, \dots, N)$
  - (b) Sample  $w_m \sim p(w | y_m, \mathbf{z}, \beta)$  from a multinomial distribution conditioned on the  $z_{y_m}$  factor.

CORR-LDA thus specifies the following joint distribution on image regions, caption words, and latent variables:

$$p(\mathbf{r}, \mathbf{w}, \theta, \mathbf{z}, \mathbf{y}) = p(\theta | \alpha) \left( \prod_{n=1}^N p(z_n | \theta) p(r_n | z_n, \mu, \sigma) \right) \cdot \left( \prod_{m=1}^M p(y_m | N) p(w_m | y_m, \mathbf{z}, \beta) \right).$$

The independence assumptions of the CORR-LDA model are a compromise between the extreme correspondence enforced by the GM-MIXTURE model, where the entire image

and caption are conditional on the same factor, and the lack of correspondence in the GM-LDA model, where the image regions and caption words can conceivably be conditional on two disparate sets of factors. Under the CORR-LDA model, the regions of the image can be conditional on any ensemble of factors but the words of the caption must be conditional on factors which are present in the image. In effect, our model captures the notion that the image is generated first and the caption annotates the image.

Finally, note that the correspondence implemented by CORR-LDA is not a one-to-one correspondence, but is more flexible: all caption words could come from a subset of the image regions, and multiple caption words can come from the same region.

## 4. INFERENCE AND ESTIMATION

In this section, we describe approximate inference and parameter estimation for the CORR-LDA model. As a side effect of the inference method, we can compute approximations to our three distributions of interest:  $p(w|\mathbf{r})$ ,  $p(w|\mathbf{r}, r_n)$ , and  $p(\mathbf{w}, \mathbf{r})$ .

### 4.1 Variational inference

Exact probabilistic inference for CORR-LDA is intractable; as before, we avail ourselves of variational inference methods [7] to approximate the posterior distribution over the latent variables given a particular image/caption.

In particular, we define the following factorized distribution on the latent variables:

$$q(\theta, \mathbf{z}, \mathbf{y}) = q(\theta|\gamma) \left( \prod_{n=1}^N q(z_n|\phi_n) \right) \left( \prod_{m=1}^M q(y_m|\lambda_m) \right),$$

with free (variational) parameters  $\gamma$ ,  $\phi$ , and  $\lambda$ . Each variational parameter is appropriate to its respective random variable. Thus  $\gamma$  is a  $K$ -dimensional Dirichlet parameter,  $\phi_n$  are  $N$   $K$ -dimensional multinomial parameters, and  $\lambda_m$  are  $M$   $N$ -dimensional multinomial parameters.

Following the general recipe for variational approximation, we minimize the KL-divergence between this factorized distribution and the true posterior thus inducing a dependence on the data  $(\mathbf{r}, \mathbf{w})$ . Taking derivatives with respect to the variational parameters, we obtain the following coordinate ascent algorithm:

1. Update the posterior Dirichlet parameters:

$$\gamma_i = \alpha_i + \sum_{n=1}^N \phi_{ni}.$$

2. For each region, update the posterior distribution over factors. Note that this update takes into account the likelihood that each caption word was generated by the same factor as the region:

$$\phi_{ni} \propto p(r_n|z_n = i, \mu, \sigma) \exp\{E_q[\log \theta_i|\gamma]\} \cdot \exp\left\{\sum_{m=1}^M \lambda_{mn} \log p(w_m|y_m = n, z_m = i, \beta)\right\},$$

where  $E_q[\log \theta_i|\gamma] = \Psi(\gamma_i) - \Psi(\sum \gamma_j)$ , and  $\Psi$  is the digamma function.

3. For each word, update the approximate posterior distribution over regions:

$$\lambda_{mn} \propto \exp\left\{\sum_{i=1}^K \phi_{ni} \log p(w_m|y_m = n, z_n = i, \beta)\right\}.$$

These update equations are invoked repeatedly until the change in KL divergence is small.

With the approximate posterior in hand, we can find a lower bound on the joint probability,  $p(\mathbf{w}, \mathbf{r})$ , and also compute the conditional distributions of interest:  $p(w|\mathbf{r})$  and  $p(w|\mathbf{r}, r_n)$ . In annotation, we approximate the conditional distribution over words as follows:

$$p(w|\mathbf{r}) \approx \sum_{n=1}^N \sum_{z_n} q(z_n|\phi_n) p(w|z_n, \beta).$$

In region labeling, the distribution over words conditioned on an image and a region is approximated by:

$$p(w|\mathbf{r}, r_n) \approx \sum_{z_n} q(z_n|\phi_n) p(w|z_n, \beta).$$

### 4.2 Parameter estimation

Given a corpus of image/caption data,  $\mathcal{D} = \{(\mathbf{r}_d, \mathbf{w}_d)\}_{d=1}^D$ , we find maximum likelihood estimates of the model parameters with a *variational EM* procedure that maximizes the lower bound on the log likelihood of the data induced by the variational approximation described above. In particular, the E-step computes the variational posterior for each image and caption given the current setting of the parameters. The M-step subsequently finds maximum likelihood estimates of the model parameters from expected sufficient statistics taken under the variational distribution. The variational EM algorithm alternates between these two steps until the bound on the expected log likelihood converges.

#### 4.2.1 Smoothing with empirical Bayes

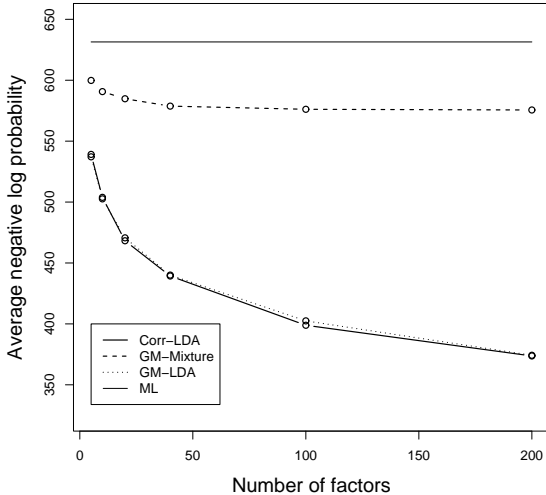
In Section 5, we show that overfitting can be a serious problem, particularly when working with the conditional distributions for image annotation. We deal with this issue by taking a more thoroughgoing Bayesian approach, imposing a prior distribution on the word multinomial parameters  $\beta$ . We represent  $\beta$  as a matrix whose columns are the  $K$  multinomial parameters for the latent factors. We treat each column as a sample from an exchangeable Dirichlet distribution  $\beta_i \sim \text{Dir}(\eta, \eta, \dots, \eta)$  where  $\eta$  is a scalar parameter.

In place of a point estimate, we now have a smooth posterior  $p(\beta|\mathcal{D})$ . A variational approach can again be used to find an approximation to this posterior distribution [1]. Introducing variational Dirichlet parameters  $\rho_i$  for each of the  $K$  multinomials, we find that the only change to our earlier algorithm is to replace the maximization with respect to  $\beta$  with the following variational update:

$$\rho_{ij} = \eta + \sum_{d=1}^D \sum_{m=1}^M 1(w_{dm} = j) \sum_{n=1}^K \phi_{ni} \lambda_{mn},$$

and we replace all instances of  $\beta$  in the variational inference algorithm by  $\exp\{E[\log \beta|\rho]\}$ .

Bayesian methods often assume a noninformative prior which, in the case of the exchangeable Dirichlet, means setting  $\eta = 1$ . With  $K$  draws from the prior distribution, however, we are in a good position to take the empirical Bayes perspective [9] and compute a maximum likelihood estimate of  $\eta$ . This amounts to a variant of the ML procedure for a Dirichlet with expected sufficient statistics under  $\rho$ . Analogous smoothing algorithms are readily derived for the GM-MIXTURE and GM-LDA models.



**Figure 4: The per-image average negative log probability of the held-out test set as a function of the number of hidden factors (lower numbers are better). The horizontal line is the model that treats the regions and captions as an independent Gaussian and multinomial, respectively.**

## 5. RESULTS

In this section, we present an evaluation of all three models on 7000 images and captions from the Corel database. We held out 25% of the data for testing purposes and used the remaining 75% to estimate parameters. Each image is segmented into 6-10 regions and is associated with 2-4 caption words. The vocabulary contains 168 unique terms.

### 5.1 Test set likelihood

To evaluate how well a model fits the data, we computed the per-image average negative log likelihood of the test set on all three models for various values of  $K$ . A model which better fits the data will assign a higher likelihood to the test set (i.e., lower numbers are better in negative likelihood).

Figure 5 illustrates the results. As expected, GM-LDA provides a much better fit than GM-MIXTURE. Furthermore, CORR-LDA provides as good a fit as GM-LDA. This is somewhat surprising since GM-LDA is a less constrained model. However, both models have the same number of parameters; their similar performance indicates that, on average, the number of hidden factors used to model a particular image is adequate to model its caption.<sup>2</sup>

### 5.2 Automatic annotation

Given a segmented image without its caption, we can use the mixture models described in Section 3 to compute a distribution over words conditioned on the image,  $p(w|\mathbf{r})$ . This distribution reflects a prediction of the missing caption words for that image.

<sup>2</sup>Empirically, when  $K = 200$ , we find that in only two images of the test set does the GM-LDA model use more hidden factors for the caption than it does for the image.

### 5.2.1 Caption perplexity

To measure the annotation quality of the models, we computed the *perplexity* of the given captions under  $p(w|\mathbf{r})$  for each image in the test set. Perplexity, which is used in the language modeling community, is equivalent algebraically to the inverse of the geometric mean per-word likelihood (again, lower numbers are better):

$$\text{perplexity} = \exp\left\{-\sum_{d=1}^D \sum_{m=1}^{M_d} \log p(w_m | \mathbf{r}_d) / \sum_{d=1}^D M_d\right\}.$$

Figure 5 (Left) shows the perplexity of the held-out captions under the maximum likelihood estimates of each model for different values of  $K$ . We see that overfitting is a serious problem in the GM-MIXTURE model, and its perplexity immediately grows off the graph (e.g., when  $K = 200$ , the perplexity is 2922). Note that in related work [2], many of the models considered are variants of GM-MIXTURE and rely heavily on an ad-hoc smoothing procedure to correct for overfitting.

Figure 5 (Right) illustrates the caption perplexity under the smoothed estimates of each model using the empirical Bayes procedure from Section 4.2.1. The overfitting of GM-MIXTURE has been corrected. Once smoothed, it performs better than GM-LDA despite the GM-LDA model’s superior performance in joint likelihood.

We found that GM-LDA does not provide good conditional distributions for two reasons. First, it is “over-smoothed.” Computing  $p(w|\mathbf{r})$  requires integrating a diffuse posterior (due to the small number of regions) over all the factor dimensions. Thus, the factors to which each region is associated are essentially washed out and, as  $K$  gets large, the model’s performance approaches the performance of the simple maximum likelihood estimate of the caption words.

Second, GM-LDA easily allows caption words to be generated by factors that did not contribute to generating the image regions (e.g., when  $K = 200$ , 54% of the caption words in the test set are assigned to factors that do not appear in their corresponding images). With this freedom, the estimated conditional Gaussian parameters do not necessarily reflect regions that are correctly annotated by the corresponding conditional multinomial parameters. While it better models the joint distribution of words and regions, it fails to model the relationship between them.

Most notably, CORR-LDA finds much better predictive distributions of words than either GM-LDA or GM-MIXTURE. It provides as flexible a joint distribution as GM-LDA but guarantees that the latent factors in the conditional Gaussian (for image regions) correspond with the latent factors in the conditional multinomial (for caption words). Furthermore, by allowing caption words to be allocated to different factors, the CORR-LDA model achieves superior performance to the GM-MIXTURE which is constrained to associating the entire image/caption to a single factor. Thus, with CORR-LDA, we can achieve a competitive fit of the joint distribution and find superior conditional distributions of words given images.

### 5.2.2 Annotation examples

Figure 6 shows ten sample annotations—the top five words from  $p(w|\mathbf{r})$ —computed by each of the three models for  $K = 200$ . These examples illustrate the limitations and power of the probabilistic models described in Section 3 when used for a practical discriminative task.

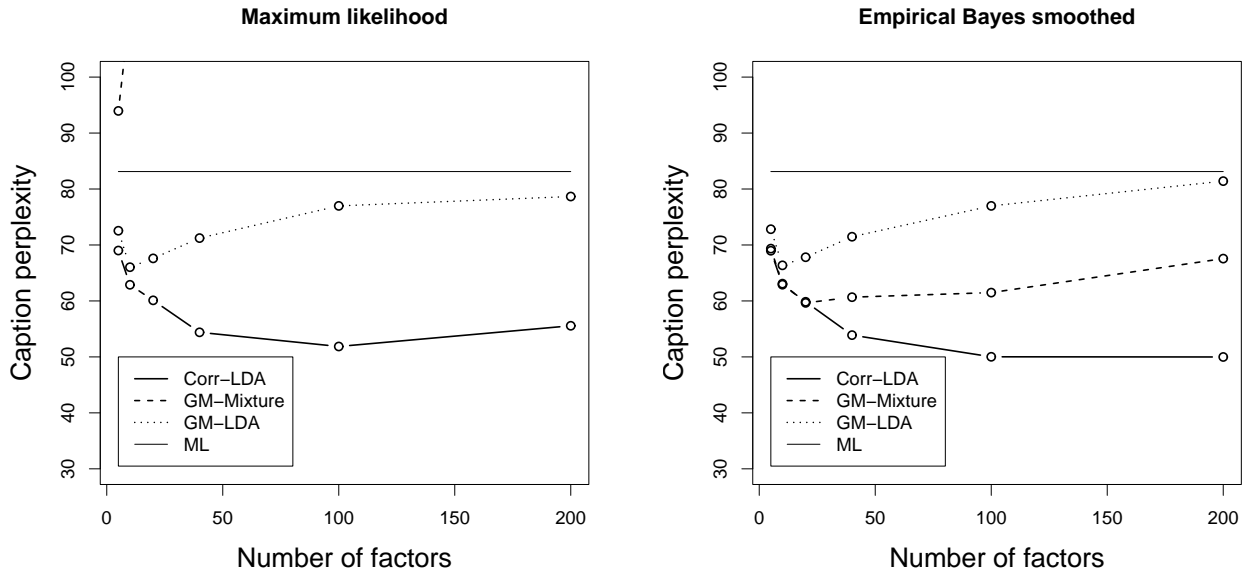
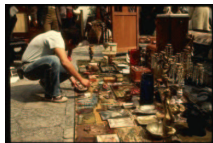


Figure 5: (Left) Caption perplexity on the test set for the ML estimates of the models (lower numbers are better). Note the serious overfitting problem in GM-Mixture (values for  $K$  greater than five are off the graph) and the slight overfitting problem in Corr-LDA. (Right) Caption perplexity for the empirical Bayes smoothed estimates of the models. The overfitting problems in GM-Mixture and Corr-LDA have been corrected.



**True caption**  
market people  
**Corr-LDA**  
people market pattern textile display  
**GM-LDA**  
people tree light sky water  
**GM-Mixture**  
people market street costume temple



**True caption**  
scotland water  
**Corr-LDA**  
scotland water flowers hills tree  
**GM-LDA**  
tree water people mountain sky  
**GM-Mixture**  
water sky clouds sunset scotland



**True caption**  
bridge sky water  
**Corr-LDA**  
sky water buildings people mountain  
**GM-LDA**  
sky water people tree buildings  
**GM-Mixture**  
sky plane jet water snow



**True caption**  
sky tree water  
**Corr-LDA**  
tree water sky people buildings  
**GM-LDA**  
sky tree fish water people  
**GM-Mixture**  
tree vegetables pumpkins water garden



**True caption**  
birds tree  
**Corr-LDA**  
birds nest leaves branch tree  
**GM-LDA**  
water birds nest tree sky  
**GM-Mixture**  
tree ocean fungus mushrooms coral



**True caption**  
fish reefs water  
**Corr-LDA**  
fish water ocean tree coral  
**GM-LDA**  
water sky vegetables tree people  
**GM-Mixture**  
fungus mushrooms tree flowers leaves



**True caption**  
mountain sky tree water  
**Corr-LDA**  
sky water tree mountain people  
**GM-LDA**  
sky tree water people buildings  
**GM-Mixture**  
buildings sky water tree people



**True caption**  
clouds jet plane  
**Corr-LDA**  
sky plane jet mountain clouds  
**GM-LDA**  
sky water people tree clouds  
**GM-Mixture**  
sky plane jet clouds pattern

Figure 6: Example images from the test set and their automatic annotations under different models.

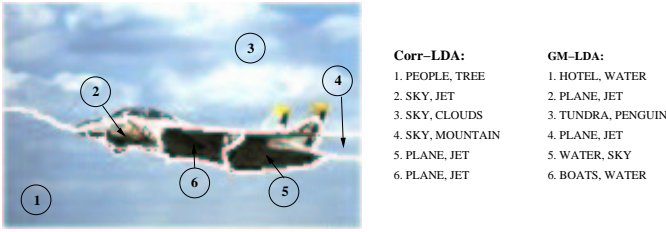


Figure 7: An example of automatic region labeling.

The GM-LDA model, as shown quantitatively in the previous section, gives the least impressive performance of the three. First, we see washing out effect described above by the fact that many of the most common words in the corpus — words like “water” and “sky” — occur in the predicted captions for all of the pictures. Second, the predicted caption rarely predicts the object or objects that are in the picture. For example, it misses “jet” in the picture captioned *clouds, jet, plane*, a word that both other models predict with high accuracy.

The GM-MIXTURE model performs better than GM-LDA, but we can see how this model relies on the average image features and fails to predict words for regions that may not generally occur in other similar images. For example, it omits “tree” from the picture captioned *scotland, water* since the trees are only a small part on the left side of the frame. Furthermore, the GM-MIXTURE predicts completely incorrect words if the average features do not easily correspond to a common theme. For example, the background of *fish, reefs, water* is not the usual blue and GM-MIXTURE predicts words like “fungus”, “tree”, and “flowers.”

Finally, as reflected by the term perplexity results above, the CORR-LDA model gives the best performance and correctly labels most of the example pictures. Unlike the GM-MIXTURE model, it can assign each region to a different cluster and the final distribution over words reflects the ensemble of clusters which were assigned to the image regions. Thus, the CORR-LDA model finds the trees in the picture labeled *scotland, water* and can correctly identify the fish, even without its usual blue background.

As described in Section 3, the CORR-LDA and GM-LDA models can furthermore compute a region-based distribution over words,  $p(w | \mathbf{r}, r_n)$ . Figure 7 illustrates a sample region labeling on an image in the test set.<sup>3</sup> Though both models hypothesize the word “plane” and “jet,” CORR-LDA places them in the reasonable regions 2, 5, and 6 while GM-LDA places them in regions 2 and 4. Furthermore, CORR-LDA recognizes the top region as “sky, clouds” while GM-LDA provides the enigmatic “tundra, penguin.”

### 5.3 Text-based image retrieval

There has been a significant amount of computer science research on *content-based image retrieval* in which a particular query image (possibly a sketch or primitive graphic) is used to find matching relevant images [5]. In another line of research, *multimedia information retrieval*, representations of different data types (such as text and images) are used to retrieve documents that contain both [8].

<sup>3</sup>We cannot quantitatively evaluate this task (i.e., compute the region perplexity) because our data does not provide ground-truth for the region labels.

Less attention, however, has been focused on *text-based image retrieval*, an arguably more difficult task where a user submits a text query to find matching images for which there is no related text. Previous approaches have essentially treated this task as a classification problem, handling specific queries from a vocabulary of about five words [10]. In contrast, by using the conditional distribution of words given an image, our approach can handle arbitrary queries from a large vocabulary.

We use a unique form of the language modeling approach to information retrieval [11] where the document language models are derived from images rather than words. For each unannotated image, we obtain an image-specific distribution over words by computing the conditional distribution  $p(w | \mathbf{r})$  which is available for the models described in Section 3. This distribution provides a description of each image in word-space which we use to find images that are similar to the words of the query.

More formally, denote an  $N$ -word query by  $\mathbf{q} = \{q_1, \dots, q_N\}$ . For each image  $\mathbf{r}_i$ , its score relative to the query is:

$$\text{score}_i = \prod_{n=1}^N p(q_n | \mathbf{r}_i),$$

where  $p(q_n | \mathbf{r}_i)$  is the probability of the  $n$ th query word under the distribution  $p(w | \mathbf{r}_i)$ . After computing the score for each image, we return a list of images ranked in descending order by conditional likelihood.

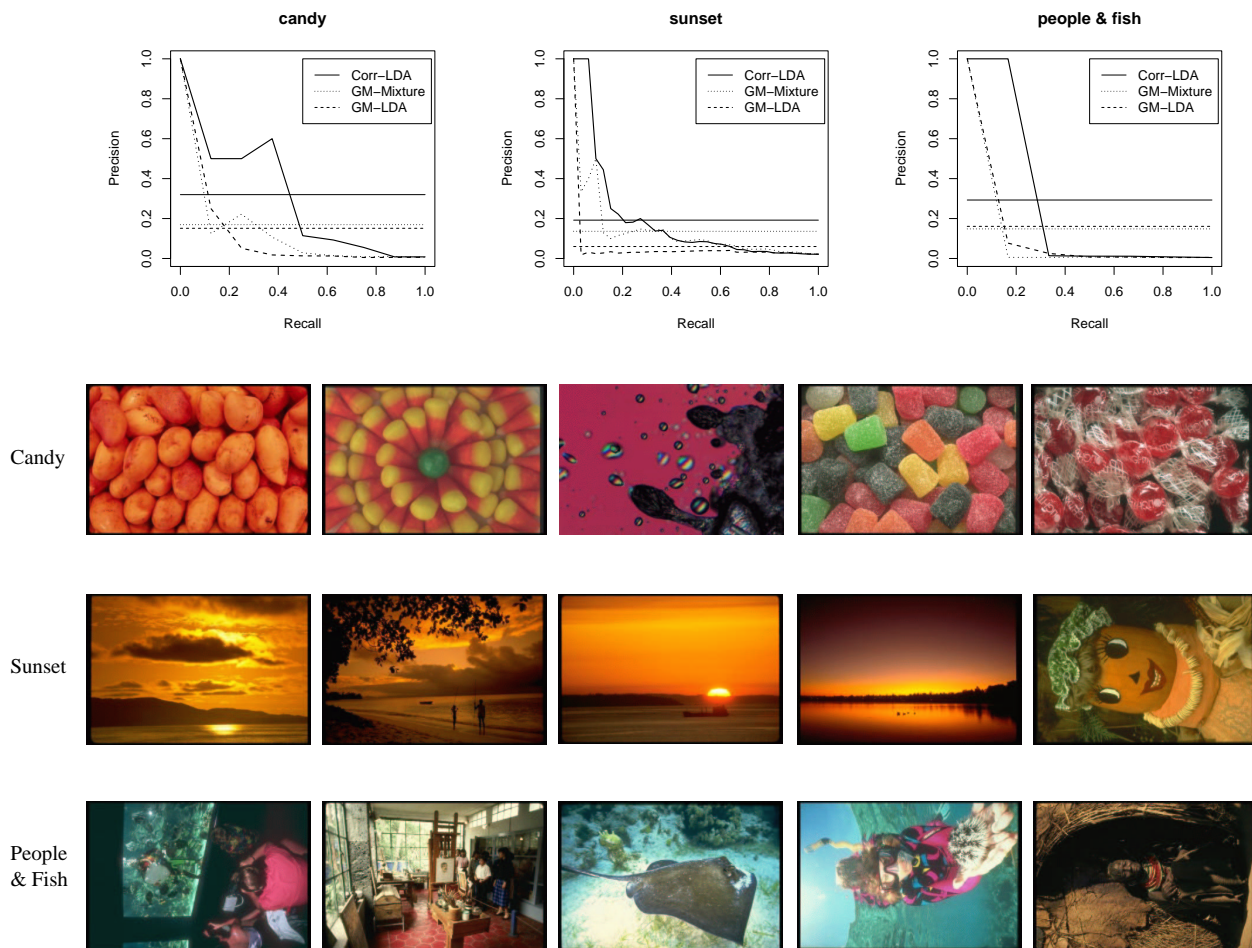
Figure 8 illustrates three queries performed on the three models with 200 factors and the held-out test set of images. We consider an image to be relevant if its true caption contains the query words (recall that we make no reference to the true caption in the retrieval process). As illustrated by the precision/recall curves, the CORR-LDA model achieves superior retrieval performance. It is particularly strong with difficult queries such as “people and fish.” In this example, there are only six relevant images in the test set and two of them appear in the top five. This is due to the ability of CORR-LDA to assign different regions to different clusters. The model can independently learn the salient features of “fish” and “people” and effectively combine them to perform retrieval.

## 6. SUMMARY

We have developed CORR-LDA, a powerful model for annotated data that combines the advantages of probabilistic clustering for dimensionality reduction with an explicit model of the conditional distribution from which data annotations are generated. In the setting of image/caption data, we have shown that this model can achieve a competitive joint likelihood and superior conditional distribution of words given an image.

CORR-LDA provides a clean probabilistic model for performing various tasks associated with multi-type data such as images and their captions. We have demonstrated its use in automatic image annotation, automatic image region annotation, and text-based image retrieval. It is important to note that this model is not specially tailored for image/caption data. Given good features, the CORR-LDA model can be applied to any kind of annotated data such as video/closed-captions, music/text, and gene/functions.





**Figure 8: Three examples of text-based image retrieval. (Top) Precision/recall curves for three queries on a 200-factor Corr-LDA model. The horizontal lines are the mean precision for each model. (Bottom) The top five returned images for the same three queries.**

## Acknowledgments

The authors thank David Forsyth and Kobus Barnard for numerous conversations and the well-curated data on which this work is based. This work was supported by a grant from the Intel Corporation, the National Science Foundation (NSF grant IIS-9988642), and the Multidisciplinary Research Program of the Department of Defense (MURI N00014-00-1-0637). David M. Blei was additionally supported by a grant from the Microsoft Corporation.

## 7. REFERENCES

- [1] H. Attias. A variational Bayesian framework for graphical models. In *Advances in Neural Information Processing Systems 12*, 2000.
- [2] K. Barnard, P. Duygulu, N. de Freitas, D. Forsyth, D. Blei, and M. Jordan. Matching words and pictures. *Journal of Machine Learning Research*, 3:1107–1135, 2003.
- [3] D. Blei, A. Ng, and M. Jordan. Latent Dirichlet allocation. *Journal of Machine Learning Research*, 3:993–1022, January 2003.
- [4] D. Cohn and T. Hofmann. The missing link—A probabilistic model of document content and hypertext connectivity. In *Advances in Neural Information Processing Systems 13*, 2001.
- [5] A. Goodrum. Image information retrieval: An overview of current research. *Informing Science*, 3(2):63–67, 2000.
- [6] J. Jeon, V. Lavrenko, and R. Manmatha. Automatic image annotation and retrieval using cross-media relevance models. In *ACM SIGIR 2003*, July 2003.
- [7] M. Jordan, Z. Ghahramani, T. Jaakkola, and L. Saul. Introduction to variational methods for graphical models. *Machine Learning*, 37:183–233, 1999.
- [8] C. Meghini, F. Sebastiani, and U. Straccia. A model of multimedia information retrieval. *Journal of the ACM (JACM)*, 48(5):909–970, 2001.
- [9] C. Morris. Parametric empirical Bayes inference: Theory and applications. *Journal of the American Statistical Association*, 78(381):47–65, 1983.
- [10] M. Naphade and T. Huang. A probabilistic framework for semantic video indexing, filtering and retrieval. *IEEE Transactions on Multimedia*, 3(1):141–151, March 2001.
- [11] J. Ponte and B. Croft. A language modeling approach to information retrieval. In *ACM SIGIR 1998*, pages 275–281.
- [12] J. Shi and J. Malik. Normalized cuts and image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 22(9):888–905, 2000.
- [13] B. Taskar, E. Segal, and D. Koller. Probabilistic clustering in relational data. In *IJCAI-01*, pages 870–876, 2001.