# Image Classification and Retrieval from User-Supplied Tags

Hamid Izadinia
Univ. of Washington

Ali Farhadi
Univ. of Washington

Aaron Hertzmann
Adobe Research

Matthew D. Hoffman
Adobe Research

## Abstract

*This paper proposes direct learning of image classification from user-supplied tags, without filtering. Each tag is supplied by the user who shared the image online. Enormous numbers of these tags are freely available online, and they give insight about the image categories important to users and to image classification. Our approach is complementary to the conventional approach of manual annotation, which is extremely costly. We analyze of the Flickr 100 Million Image dataset, making several useful observations about the statistics of these tags. We introduce a large-scale robust classification algorithm, in order to handle the inherent noise in these tags, and a calibration procedure to better predict objective annotations. We show that freely available, user-supplied tags can obtain similar or superior results to large databases of costly manual annotations.*

## 1. Introduction

Object recognition has made dramatic strides in the past few years. This progress is partly due to the creation of large-scale hand-labeled datasets. Collecting these datasets involves listing object categories, searching the web for images of each category, pruning irrelevant images and providing detailed labels for each image. There are several major issues with this approach. First, gathering high-quality annotations for large datasets requires substantial effort and expense. Second, it remains unclear how best to determine the list of categories. Existing datasets comprise only a fraction of recognizable visual concepts, and often miss concepts that are important to end-users. These datasets draw rigid distinctions between different types of concepts (e.g., scenes, attributes, objects) that exclude many important concepts.

This paper introduces an approach to learning about visual concepts by employing user-supplied tags. That is, we directly use the tags provided by the users that uploaded the images to photo-sharing services, without any subsequent manual filtering or curation. Tags in the photosharing services reflect the image categories that are important to users and include scenes (beach), objects (car), attributes

(rustic), activities (wedding), and visual styles (portrait), as well as concepts that are harder to categorize (family). Online sharing is growing and many services host content other than photographs (e.g., Behance, Imgur, Shapeways). The tags in these services are abundant, and learning about them could benefit a broad range of consumer applications such as tag suggestion and search-by-tag.

User-supplied tags are freeform and using them presents significant challenges. These tags are entirely uncurated, so users provide tags for their images in different ways. Different users provide different numbers of tags per image, and, conversely, choose different subsets of tags. One tag may have multiple meanings, and, conversely, multiple terms may be used for the same concept. Most sharing sites provide no quality control whatsoever for their tags. Hence, it is important to design learning algorithms robust to these factors.

**Contributions.** In addition to introducing the direct use of user-supplied tags, this paper presents several contributions. First, we analyze statistics of tags in a large Flickr dataset, making useful observations about how tags are used and when they are reliable. Second, we introduce a robust logistic regression method for classification with user-supplied tags, which is robust to randomly omitted positive labels. Since tag noise is different for different tags, the tag outlier probabilities are learned simultaneously with the classifier weights. Third, we describe calibration of the trained model probabilities from a small annotation set.

We demonstrate results for several tags: predicting the tags that a user would give to an image, predicting objective annotations for an image, and retrieving images for a tag query. For the latter two tasks, which require objective anotations, we calibrate and test on the manually-annotated NUS-WIDE [5] dataset.

We show that training on a large collection of freely available, user-supplied tags alone obtains comparable performance to using a smaller, manually-annotated training set. That is, we can learn to predict thousands of tags *without any curated annotations at all.* Moreover, if we calibrate the model with a small annotated dataset, we can obtain superior performance to conventional annotations at

a tiny fraction (1/200) of the labeling cost. Our methods could support several annotation applications, such as auto-suggesting tags to users, clustering user photos by activity or event, and photo database search. We also demonstrate that using robust classification substantially improves image retrieval performance with multi-tag queries.

## 2. Related Work

The amazing progress of the recent years of vision has been driven in part by datasets. These datasets are built through a combination of webscraping and crowd-sourcing, with the aim of labeling the data as cleanly as possible. ImageNet [21] is the most prominent whole-image classification dataset, but other recent examples include NUS-WIDE [5], the SUN scene attribute database [19, 26], and PLACES [27]. The curation process has a number of drawbacks, such as the cost of gathering clean labels and the difficulty in determining a useful space of labels. It is unclear that this procedure alone will scale to the space of all important concepts for vision [21]. We take a complementary approach of using a massive database of freely available images with noisy, unfiltered tags.

Merging noisy labels is a classic problem in item-response theory, and has been applied in the crowdsourcing literature [20, 25]. We extend robust logistic regression [20] to large-scale learning with Stochastic EM. In image recognition, a related problem occurs when harvesting noisy data from the web [3, 4, 7, 15, 24]; these methods take complementary approaches to ours, and focus on object and scene categories.

To our knowledge, no previous work directly learns image classifiers from raw Flickr tags without curation. Most similar to our own work, Zhu et al. [28] use matrix factorization to clean up a collection of tags. In principle, this method could be used as a first step toward learning classifiers, though it has not been tested as such. This method requires batch computation and is unlikely to be practical for large numbers of tags and images. Gong et al. [9] use raw Flickr tags as side-information for associating images with descriptive text.

Most previous work has focused on names and attributes for objects and scenes, including previous work on image tagging (e.g., [6, 8, 11, 19, 21, 26, 27]). Unfortunately these datasets are disjoint and little attention has been paid to the list of objects, scenes, and attributes. Our solution is to learn what users care about, using a robust loss function that takes into account the noise in the labels. We learn many other kinds of tags, such as tags for events, activities, and image style. There have been a few efforts aimed at modeling a few kinds of image style and aesthetics [13, 17].
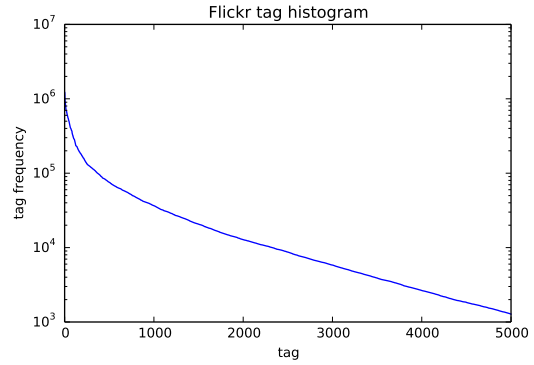


Figure 1: Tag histogram for the most popular tags, excluding non-image tags. The distribution is heavy-tailed, and there are 5400 tags with more than 1000 images each.

## 3. Analysis of User-Supplied Tags

When can user-supplied tags be useful, and when can they be trusted? In this section, we analyze the tags provided on Flickr, and compare them to two datasets with ground truth labels. Some of these observations motivate our algorithm in Section 4, and others provide fodder for future research.

**Flickr 100 Million (F100M).** Our main dataset is the Yahoo/Flickr Creative Commons 100M dataset[1]. This dataset comprises 99.3 million images, each of which includes a list of the tags supplied by the user that uploaded the image.

### 3.1. Types of tags

The F100M dataset provides an enormous number of images and tags (Figure 1) that could be used for learning. Some of the most frequent tags are shown in Table 1. There are 5400 tags that occur in at least 1000 images. The set of tags provides a window into the image concepts that are important to users. Many of these represent types of image label that are not represented in previous datasets.

Some of the most important tag types are as follows: **events and activities** such as travel, music, party, festival, football, school; **specific locations** such as california and italy; **scene types** such as nature, part, urban, sunset, etc.; **the seasons** (fall, winter, summer, spring); **image style** such as portrait, macro, vintage, hdr; and **art and culture** such as painting, drawing, graffiti, fashion, punk. Many frequent tags also represent categories that do not seem learnable from image data alone, which we call **non-image tags**, including years (2011, 2012, ...), and specific camera and imaging platforms (nikon, iphone, slr).

---

[1]http://yahoolabs.tumblr.com/post/89783581601

| Flickr tag | # Flickr | synset | # node | # subtree |
|---|---|---|---|---|
| travel | 1221148 | *travel.n.01* | 0 | 0 |
| wedding | 734438 | *wedding.n.03* | 1257 | 1257 |
| flower | 907773 | *flower.n.01* | 1924 | 339376 |
| art | 902043 | *art.n.01* | 0 | 11353 |
| music | 826692 | *music.n.01* | 0 | 0 |
| party | 669065 | *party.n.01* | 0* | 0 |
| nature | 872029 | *nature.n.01* | 0 | 0 |
| beach | 768752 | *beach.n.01* | 1713 | 1773 |
| city | 701823 | *city.n.01* | 1224 | 1224 |
| tree | 697009 | *tree.n.01* | 1181 | 563038 |
| vacation | 694523 | *vacation.n.01* | 0 | 0 |
| park | 686458 | *park.n.01* | 0 | 0 |
| people | 641571 | *people.n.01* | 1431 | 1431 |
| water | 640259 | *water.n.06* | 759 | 7585 |
| architecture | 616299 | *architecture.n.01* | 1298 | 1298 |
| car | 610114 | *car.n.01* | 1307 | 40970 |
| festival | 609638 | *festival.n.01* | 0 | 0 |
| concert | 605163 | *concert.n.01* | 1322 | 1322 |
| summer | 601816 | *summer.n.01* | 0 | 0 |
| sport | 564703 | *sport.n.01* | 1888 | 200402 |

Table 1: The 20 most frequent tags in F100M, after merging plurals and omitting non-image/location tags. Corresponding ImageNet synsets are given, along with synset node and subtree counts. These statistics are typical: we estimate that nearly half of popular Flickr tags are absent from ImageNet. Moreover, even when there is correspondence, some ImageNet tags do not capture all meanings of a term (Section 3.2). Some of these tags are covered by scene attribute databases [26, 19, 27]. (*There are 66 party images in ImageNet, in the wrong synset *party.n.04*.)

## 3.2. Correspondence with ImageNet

A main motivation for using F100M is that it contains information missing from existing, curated datasets. Does it? We compare F100M to the ImageNet image classification dataset [21], which comprises 14 million images gathered from Flickr, labeled according to the WordNet hierarchy [16] through a carefully-designed crowdsourcing procedure.

In order to quantify the dataset gap, we studied the 100 most frequent tags in F100M (after omitting the non-image and location tags described above). For each tag, we manually determined a correspondence to WordNet, as follows. In WordNet, each concept is represented by a synonym set, or *synset*. WordNet synsets are ordered, and most tags (78%) correspond to the first WordNet noun synset for that word. For example, the tag beach corresponds to the synset *beach.n.01*. In other cases, we corrected the match manually. The most-frequent examples are shown in Table 1, and more are shown in the Appendix. *Based on this analysis and some simple calculations, we estimate that about half of the common Flickr non-image tags are absent from ImageNet.* Details of how this estimate was formed are given in the Appendix. Some of these missing tags are covered by scene [19, 26, 27] and style databases [13, 17].

Even when there is a corresponding tag in ImageNet, the tag may be poorly represented. There are 11k images in the ImageNet *art.n.01* hierarchy, but there are only 8 subtrees of *art.n.01* with at least 1000 images; the biggest ones are "olympian zeus," "cinquefoil," and "finger-painting;" and there are no subtrees for "painting," "drawing," or "illustration." The ImageNet synset for "band" includes only images for "marching bands" and not, say, "rock bands."

Many image categories that are significant to users—for example, in analyzing personal photo collections—are not well represented in the ImageNet categories. Examples include family, travel, festival, and summer.

Some common tags in Flickr do not even exist in the WordNet hierarchy, such as cosplay (a popular form of costume play), macro (as in macro photography), and vintage (in the sense of "retro" or "old-style"). We also observed problems in the full ImageNet database, where large sets of images are assigned to the wrong synset, such as "party," "landscape," and "tree/tree diagram."

This is not in any way meant to disparage the substantial, important efforts of the ImageNet team, but to emphasize the enormous difficulty in trying to precisely curate a dataset including all important visual concepts.

## 3.3. Label noise and ambiguities

A fundamental challenge in dealing with user-supplied tags is that the mapping from observed tags to underlying concepts is ambiguous. Here we discuss many types of these ambiguities that we have observed.

Many terms have multiple or overlapping meanings. The simplest case is for plurals, e.g., car and cars, which have different meanings but which seem to be more or less interchangeable tags on Flickr. Some tags have multiple distinct meanings [22], e.g., rock can mean both "rock-and-roll music," and "rocky landscapes." Trickier cases include terms like music, concert, and performance, which often overlap, but often do not. Some words are used nearly interchangeably, such as cat and kitten, even though their meanings are not the same. It seems that nearly all common tags exhibit some multiple meanings, though often one sense dominates the others. Synonyms are also common, e.g., cat and gato, as well as misspellings.

Multi-word tags often occur split up, e.g., images in New York are frequently tagged as New and York rather than New York. For this reason, tags like New and San are largely meaningless on their own. Merging these split tags (especially using cues from the other image metadata) is a natural problem for future research.

### 3.4. Analysis with Ground Truth

In this section, we perform analysis using the annotated subset of the **NUS-WIDE dataset** [5]. This is a set of 269,642 Flickr images with annotations with both user-supplied tags, and "ground truth" annotations by undergraduate and high school students according to 81 concepts. There are a number of potential sources of noise with this dataset. Since the dataset was constructed by keyword searches, it is not an unbiased sample of Flickr, e.g., only one image in the dataset has zero keywords. Annotators were not asked to judge every image for every concept; a query expansion strategy was used to reduce annotator effort. Annotators were also asked to judge whether concepts were present in images in ways that may differ from how the images were originally tagged.

**Tagging likelihoods.** We now quantify the accuracy of Flickr tags. We consider the Flickr images in NUS-WIDE that contain manual annotations, and we treat these 81 labels as ground truth, thus expanding on the discussion in [5]. We assume an identity mapping between tags and annotations, i.e., the Flickr tag cat corresponds to the NUS-WIDE annotation cat.

Overall, given that a tag correctly applies to an image, there is empirically a 38% chance that the uploader will actually supply it. This probability varies considerably for different tags, ranging from 2% for person to 94% for cat. Frequently-omitted tags are often non-entry-level categories [18] (e.g., person) or they are not an important subject in the scene [1] (e.g., clouds, buildings). Given that a tag does not apply, there is a 1% chance that the uploader supplies it anyway. Across the NUS-WIDE tags, this probability ranges from 2% (for street) to 0.04% (for toy).

Despite these percentages, false tags and true tags are almost equally likely, since only a few of the 81 tags correctly apply to each image. Each image has an average of 1.3 tags (of the 81), and *an observed tag has only a 62% chance of being true*. This percentage varies across different tags.

None of these numbers should be taken as exact, because the NUS annotations are far from perfect (see Appendix). Additionally, many "false" tags are due to differences in word senses between Flickr and NUS-WIDE. For example, many earthquake images are clearly the result of earthquakes, but are labeled as negatives in NUS-WIDE. Many cat images that are annotated as non-cat are images of tigers, lions, and cat costumes. Many nighttime images were probably taken at night but indoors.

**Tag index effects on accuracy.** Flickr tags are provided in an ordered list. We observed that tags earlier in the list are often more accurate than later tags, and we again treat the NUS-WIDE annotations as ground truth in order to quantify
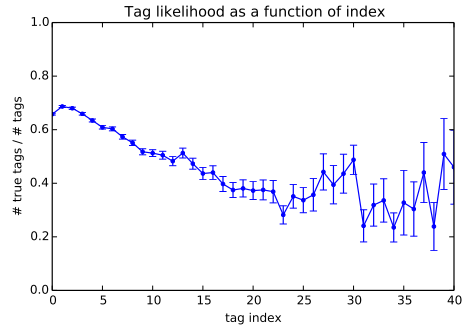


Figure 2: Empirically, tags that occur earlier in the list of an image's tags are more likely to be accurate. This plot is computed from the NUS-WIDE dataset. (Error bars show standard error.)

this.

We find that the effect is substantial, as shown in Figure 2. A tag that appears first or second in the list of tags has about 65% chance of being accurate. A tag that occurs in position 20 or later has about a 35% chance of being accurate. The scales and shape of these plots also vary considerably across different tags.

**Effect of total number of tags.** We also hypothesized that tag reliability could depend on the total number of tags provided for an image. This was motivated by our observation of commercially-oriented sharing sites, where uploaders are incentivized to include extraneous tags in order to boost search results. However, we did not find any significant effects in the Flickr data.

## 4. Robust Tag Classification

We now describe a robust classification algorithm, designed to address the following observations from the previous section: user-supplied tags often omit relevant tags, and these probabilities are different for each tag. A conventional robust loss (e.g., Huber, Geman-McClure) would not be appropriate because of the need to set the loss function's parameters individually for each tag. The method is based on previous robust logistic regression methods [20]. Previous approaches used batch computation, which cannot realistically be applied to millions of images; we adapt these methods to the large-scale setting using Stochastic EM [2].

The classifier takes as input image features $\mathbf{x}$, and predicts class labels $y \in \{0, 1\}$. We perform prediction for each possible tag independently, and so we consider simple binary classification in this paper. As image features $\mathbf{x}$, we use the output of the last fully-connected layer of Krizhevsky's ImageNet Convolutional Neural Network [14]; fc7 in the Caffe implementation [12]. We do not fine-tune the network parameters in this paper.

## 4.1. Logistic Regression

As our approach is based on logistic regression, we begin by briefly reviewing a conventional binary logistic regression classifier. The logistic regression model assumes that the probability of a positive tag (i.e., the probability that $y = 1$) given input features $\mathbf{x}$ is a linear function $\mathbf{w}^T\mathbf{x}$ passed through a sigmoid:

$$\sigma(s) \equiv 1/(1 + e^{-s}); \quad P(y = 1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T\mathbf{x}) \quad (1)$$

The loss function $L(\mathbf{w})$ for a label training set $\{(\mathbf{x}_i, y_i)\}$ is the negative log-likelihood of the data:

$$L(\mathbf{w}) = -\ln P(y_{1:N}|\mathbf{x}_{1:N}, \mathbf{w}) \quad (2)$$
$$= \sum_i(-y_i \ln \sigma(\mathbf{w}^T\mathbf{x}_i) - (1 - y_i)\ln(1 - \sigma(\mathbf{w}^T\mathbf{x}_i)))$$

Training entails optimizing $L$ with respect to $\mathbf{w}$, using stochastic gradient descent. Prediction entails computing the label probability $P(y|\mathbf{x}, \mathbf{w})$ for a new image.

## 4.2. Robust model

As discussed in Section 3, user-supplied tags are often noisy. However, the logistic regression model assumes that the observed labels $\{y_i\}$ are mostly reliable—that is, it assumes that $y_i = 1$ almost always when $\mathbf{w}^T\mathbf{x}_i$ is large.

To cope with this issue, we relax the assumption that the observed training label $y$ is the true class label. We introduce a hidden variable $z \in \{0, 1\}$ representing the true (hidden) class label. We also add a variable $\pi$ to represent the probability that a true label is added as a tag. The model parameters are then $\theta = \{\mathbf{w}, \pi\}$, and the model is:

$$P(z = 1|\mathbf{x}, \mathbf{w}) = \sigma(\mathbf{w}^T\mathbf{x}) \quad (3)$$
$$P(y = 1|z = 1, \pi) = \pi; \quad P(y = 0|z = 0) = 1 \quad (4)$$

and thus:

$$P(y = 1|\mathbf{x}, \pi, \mathbf{w}) = \pi\sigma(\mathbf{w}^T\mathbf{x}) \quad (5)$$

The loss function for training is again the negative log-likelihood of the data:

$$L(\mathbf{w}, \pi) = \sum_i \left(-y_i \ln \pi\sigma(\mathbf{w}^T\mathbf{x}_i) \right. \quad (6)$$
$$\left. -(1 - y_i)\ln(1 - \pi\sigma(\mathbf{w}^T\mathbf{x}_i))\right)$$

We also experimented with a model in which false tags are occasionally added: $P(y = 0|z = 0) = \gamma$, where $\gamma$ is another learned parameter. We found that this model did not improve performance, and so, for clarity, we omit $\gamma$ from the rest of the paper. The $\gamma$ parameter may be useful for other datasets where users produce more spurious tags. Detailed derivations of the model and gradients are straightforward, and are given in the Appendix (with $\gamma$).

Although the loss function is unchanged for positive labels ($y = 1$), the model is robust to outliers for negative
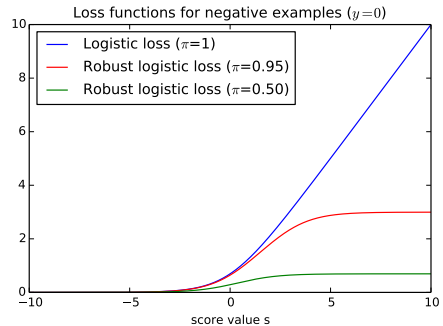


Figure 3: Loss functions for negative examples ($y = 0$). Many Flickr users omit relevant tags, which is steeply penalized by the conventional logistic loss $(\ln(1 - \sigma(s)))$. The robust logistic loss $(\ln(1 - \pi\sigma(s)))$, is tolerant to missing labels.

examples ($y = 0$); see Figure 3 for loss function plots. The classical logistic loss is unbounded, meaning that an overly confident prediction may be heavily penalized. With a true positive rate of $\pi = 0.95$, the loss is bounded above by $-\ln(1 - \pi) \approx 3$, since no matter what the image there is always at least a probability of $1 - \pi$ of a negative label. The impact of $\pi$ becomes smaller as the score $s = \mathbf{w}^T\mathbf{x}$ becomes small, since if $s \ll 0$ then $P(z = 0) \approx 1$ and $\pi$ is only relevant when $z = 1$. When the true positive rate is lower (e.g., $\pi = 0.5$ as in Figure 3), the dynamic range of the loss function is further compressed.

## 4.3. Stochastic EM algorithm

Learning the model for a given tag entails minimization of the loss with respect to $\mathbf{w}$ and $\pi$. Stochastic gradient descent could be used for all parameters, and we provide gradients in the Appendix. However, we use Stochastic Expectation-Maximization (EM) [2], since the steps are simpler to interpret and implement, and the updates to $\pi$ are numerically stable by design. All derivations and detailed versions of these equations are given in the Appendix.

Our stochastic EM algorithm applies the following steps to each minibatch:

1. For each image $i$ in the minibatch, the conditional probability of the true label $z_i$ is computed as:

$$\alpha_i \leftarrow P(z = 1|y_i, \mathbf{x}_i, \mathbf{w}, \pi) \quad (7)$$
$$= \begin{cases} 1 & y_i = 1 \\ \frac{(1-\pi)\sigma(\mathbf{w}^T\mathbf{x}_i)}{1 - \pi\sigma(\mathbf{w}^T\mathbf{x}_i)} & y_i = 0 \end{cases} \quad (8)$$

2. We define the sufficient statistics for the minibatch as

$$S_\alpha^{\text{mb}} \equiv \sum_i \alpha_i/N; \quad S_{y\alpha}^{\text{mb}} \equiv \sum_i y_i\alpha_i/N, \quad (9)$$

where $N$ is the number of datapoints in the summation. Estimates of the average sufficient statistics for the full

dataset are updated with a step size $\eta$:

$$S^{\mathrm{ds}} \leftarrow (1-\eta)S^{\mathrm{ds}} + \eta S^{\mathrm{mb}} \qquad (10)$$

In our experiments, we initialized $S^{\mathrm{ds}}_{\alpha}$ and $S^{\mathrm{ds}}_{y\alpha}$ to 1 and used a fixed step size of $\eta = 0.01$.

3. $\pi$ is computed from the current estimate of the sufficient statistics, so that $\pi$ is an estimate of the percentage of true labels that were actually supplied as tags:

$$\pi \leftarrow S^{\mathrm{ds}}_{y\alpha}/S^{\mathrm{ds}}_{\alpha} \qquad (11)$$

4. The weights $\mathbf{w}$ are updated using stochastic gradient on $L$. It is straightforward to verify that the gradient w.r.t. $\mathbf{w}$ is

$$\frac{dL}{d\mathbf{w}} = \sum_i (\sigma(\mathbf{w}^T \mathbf{x}_i) - \alpha_i)\, \mathbf{x}_i. \qquad (12)$$

## 4.4. Calibration

In many cases, we would like to predict the true class probabilities $P(z|\mathbf{x})$. Well-calibrated estimates of these probabilities are particularly useful in applications where it is important to weight the importance of multiple tags for an image, such as when trying to retrieve images characterized by multiple tags or when choosing a small number of tags to apply to an image [23].

In theory, the robust model above could learn a well-calibrated estimate of $P(z|\mathbf{x})$. However, the model still makes strong simplifying assumptions—for example, it assumes linear decision boundaries, and that label noise is independent of image content. To the extent that these assumptions are unrealistic, the model may benefit from an additional calibration step.

We tried to apply the calibration method from [23], but found that it degraded the logistic regression model's performance. This may be because it is designed to address miscalibration due to using non-probabilistic classifiers such as SVMs, rather than due to label noise.

Instead, we propose the following strategy. First, a model is learned with the large dataset. The weight vector $\mathbf{w}$ is then held fixed for each tag. However, an intercept $\beta$ is added to the model, so that the new class probability is

$$P(z = 1|\mathbf{x}, \mathbf{w}, \beta) = \sigma(\mathbf{w}^T \mathbf{x} + \beta) \qquad (13)$$

The intercept allows the model to adjust the prior probability of each class in the new dataset. Then, we continue training the model on a small curated dataset (treating it as ground truth $z$), but only update the $\beta$ parameters. Very little curated data is necessary for this process, since only one new parameter is being estimated per tag.

In our experiments, we simulate this procedure by training on the F100M data and calibrating on a subset of NUS-WIDE annotations. More general domain adaptation methods (e.g., [10]) could also be used.
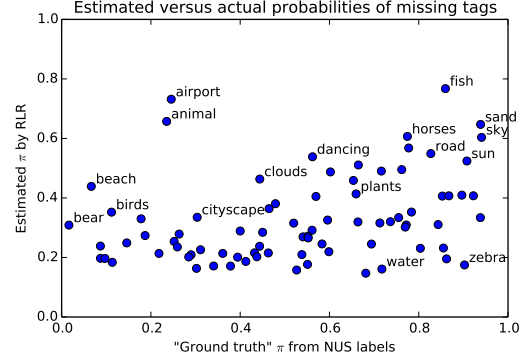


Figure 4: Tagging likelihoods $\pi$ estimated from Flickr images with RLR, versus estimation from the "ground truth" NUS annotations. The likelihoods are correlated ($r = 0.34$), though the tagging likelihood is mostly underestimated, probably due to inaccuracies in both the predictor and the annotations.

|  | Recall | Precision | F-score |
|---|---|---|---|
| LR | 9.7 | 7.9 | 8.7 |
| RLR | 11.7 | 8.0 | 9.5 |

Table 2: Flickr tag prediction results. Robust logistic regression improves over logistic regression's ability to predict which tags a user is likely to apply to an image.

|  | Recall | Precision | F-score |
|---|---|---|---|
| CNN+WARP [8] | 52.0 | 22.3 | 31.2 |
| NUS, LR | 58.2 | 26.1 | 36.0 |
| F100M, LR | 58.4 | 21.7 | 31.6 |
| F100M, RLR | 58.0 | 22.3 | 32.3 |
| F100M, LR, Calib | 42.5 | 32.2 | 36.6 |
| F100M, RLR, Calib | 44.2 | 31.3 | **36.7** |

Table 3: Image annotation results, illustrating how the freely-available user-supplied tags can augment or supplant costly manual annotations. Testing is performed on the NUS-WIDE test set. The first two rows show training only on the NUS-WIDE training set with logistic regression, and the previously-reported state-of-the-art [8]. Each of the remaining rows is trained on F100M with either LR or Robust LR. The third and fourth rows are also calibrated on the NUS test set. All scores are predictions-at-5.

## 5. Experiments

We now describe experiments to test models learned from F100M on several tasks, including tag prediction, image annotation, and image retrieval with one or more tags.

All training is performed using Caffe [12], running for 20,000 minibatches, with minibatch size of 500 images. Training is performed on a GeForce GTX780 GPU. Each minibatch takes 2 seconds, and a complete run takes 11
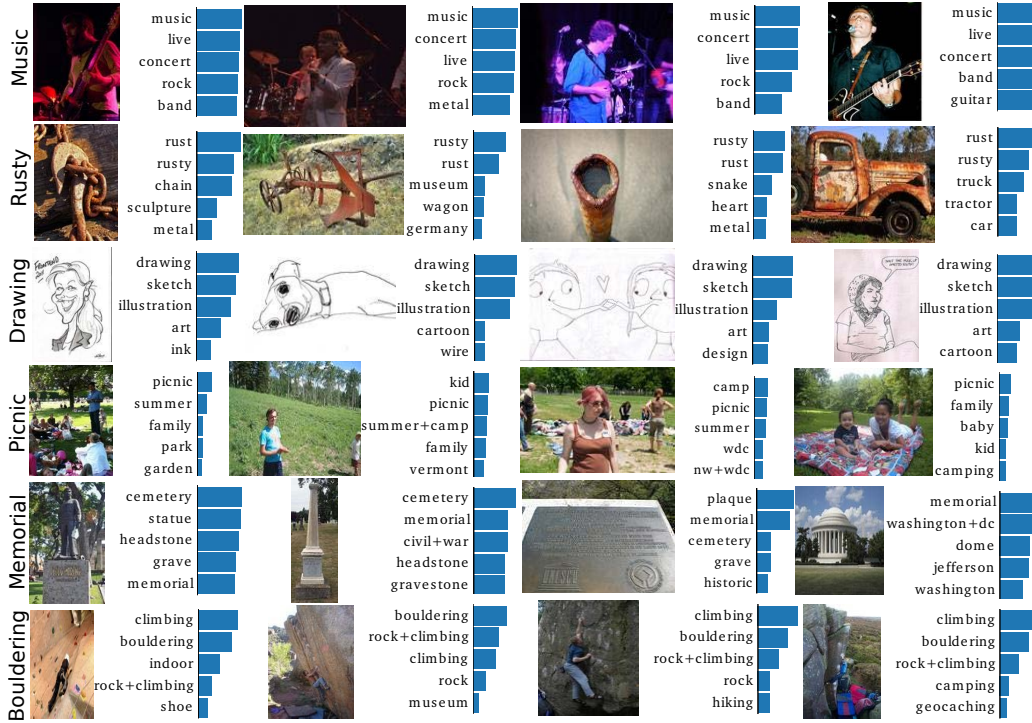
Figure 5: Single-tag retrieval results, and automatically-generated annotations. None of the query tags are in NUS-WIDE, and most (music, rusty, drawing, bouldering) are also absent from ImageNet. Many of the annotations are also absent from the other datasets as well.

hours. Based on the observations in Section 3.4, we only keep the first 20 tags in all Flickr images in our experiments. We use a subset of 4,768,700 images from F100M as training set and hold out another 200,000 for testing. The sets are split by user ID in order to ensure that images from the same user do not occur in both sets. Plural and singular tags are combined using WordNet's lemmatization.

## 5.1. Tag prediction

We first test the following prediction task: given a new image, what tags would a user be likely to apply to this image? This task could be useful for consumer applications, for example, auto-suggesting tags for users when sharing their images. Note that this task is different from ground-truth prediction; we want to suggest tags that are both objectively accurate and likely to be applied by a user.

We trained a logistic regression baseline and a robust logistic regression model on our 4.7M-image F100M training set, and evaluated the models' ability to annotate images in the 200K-image F100M test set.

For each test image, the model predicts the probability of each tag occurring: $P(y = 1|\mathbf{x}, \mathbf{w}, \pi)$. (Note that for robust logistic regression, this is Equation 5, since we want to predict tagging behavior $y$, not ground truth $z$.) The final annotations were produced by selecting the top 5 most

likely tags for each image.

We evaluate overall precision and recall at 5 for each image, averaged over all images. We also compute the F-score, which is the harmonic mean of the average precision and recall. Table 2 summarizes the results. RLR achieves higher recall without sacrificing precision. Figure 5 shows some qualitative results of calibrated RLR's ability to predict tags for images from the test set. Figure 4 compares RLR's estimated values of $\pi$ for each tag, versus the NUS annotations estimated in Section 3.4. RLR's estimates are correlated with the NUS ground truth, but discrepancies are common.

## 5.2. Image annotation

We next test the task: given an image, which labels objectively apply to it? We use the same F100M training set as above, but evaluate on the 81 labels in the manually annotated NUS-WIDE dataset, treating the NUS-WIDE annotations as ground truth. We also compare to models trained on NUS-WIDE.

We evaluate per-tag precision and recall averaged over tags. For a tag $j$, per-tag precision is defined as $N_j^c/N_j^p$ and per-tag recall is defined as $N_j^c/N_j^g$, where $N_j^p$ is the number of images that the system annotated with tag $j$, $N_j^c$ is the number of images that a user annotated with tag $j$ that the system also annotated with tag $j$, and $N_j^g$ is the
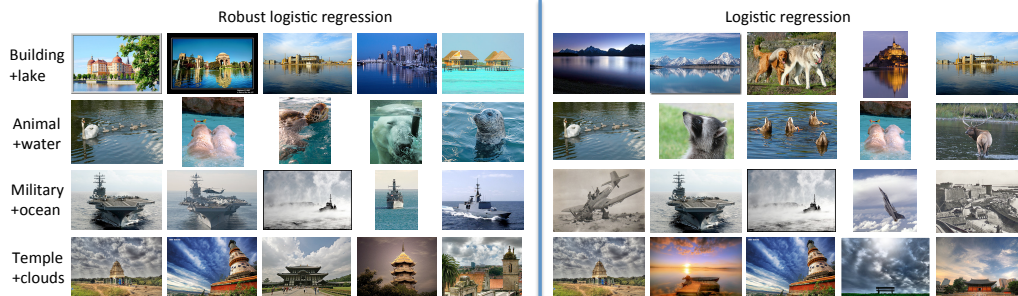
Figure 6: Multi-tag retrieval queries where Robust LR gives notably superior results to LR. Retrieval results are sorted from left-to-right.
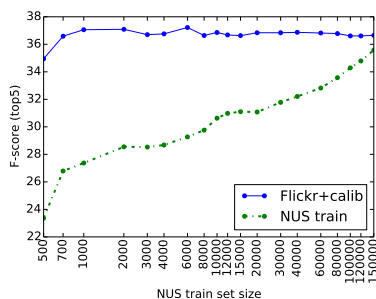


Figure 7: Effect of calibration set size on image annotation score. Training on user-supplied tags and then calibrating on a small subset of manual annotations can outperform the costly process of obtaining many manual annotations: **the annotation cost can be reduced by a factor of 200, while obtaining the same results.**

|  | 1 Tag | 2 Tags | 3 Tags |
|---|---|---|---|
| NUS, LR | *81* | *17.9* | *9.1* |
| F100M, LR | 70.1 | 8.5 | 2.3 |
| F100M, RLR | **71.9** | 9.2 | 2.7 |
| F100M, LR, Calib | 70.1 | 10.3 | 3.6 |
| F100M, RLR, Calib | **71.9** | **11** | **3.9** |

Table 4: Image retrieval results, showing precision at 5 for multi-tag retrieval. Testing is performed on the NUS-WIDE test set. Columns show performance for each method for the number of tags that need to be matched. See the caption to Table 3 for an explanation of the rows. Robust LR consistently outperforms LR, and calibration consistently improves results. These trends are clearer for longer (and therefore more difficult) queries.

number of images in the test set that a user annotated with tag $j$. Per-tag precision is undefined if a tag is never used by the system; when this happens we define the precision to be 0. We also computed the per-tag F-score. To predict annotations with RLR, we predict $z$, not $y$ (Equation 3 or 13). Scores are reported in Table 3.

Testing and training LR on NUS data produces somewhat better scores than training on F100M alone; it also produces better scores than the reported state-of-the-art on this dataset [8]. We get the best scores by training on F100M and then calibrating on the NUS training set (Section 4.4).

It is important to consider the cost of annotated labels. The user-supplied tags in F100M are basically free, whereas obtaining manual annotations is a very costly process. We compare training on a subset of NUS training annotations, versus F100M training plus calibration with the same NUS subset. As shown in Figure 7, *the calibration process can yield scores superior to training on the full annotation set, but with a 200x reduction in annotation cost.*

## 5.3. Image retrieval

Finally, we consider the tag-based image retrieval task: given a set of query tags, find images that match all the tags. We measure performance using normalized precision at 5; each system returns a set of 5 images, and its score for a given query is the number of those images that are characterized by all tags divided by the smaller of 5 and the total number of relevant images in the test set. We use the NUS-WIDE annotations as ground truth. We tested the same models from the previous section. We tested each method with queries consisting of every combination of one, two, and three tags that had at least one relevant image in the test set. Scores are shown in Table 4.

All models perform well on single-tag queries, but the differences in precision grow rapidly as the number of tags that the retrieved images must match increases. RLR consistently outperforms LR, and calibration significantly improves the models trained on Flickr. Figure 6 shows some queries for which RLR outperforms LR.

The model trained on NUS-WIDE achieves the best score. However, there are many thousands of tags for which no annotations are available, and these results show that

good results can be obtained on these tags as well.

## 6. Discussion and Future Work

Online user-supplied tags represent a great, untapped natural resource. We show that, despite their noise, these tags can be useful, either on their own or in concert with a small amount of calibration data. Though we have tested the Flickr dataset, there are numerous other online datasets with different kinds of user-supplied tags that can also be leveraged and explored for different applications. As noted in Section 3, there is a great deal of structure in these tags that could be exploited in future work. Our work could be combined with methods that model the relationships between tags, as well as improved CNN models and fine-tuning. These tags could also provide mid-level features for other classification tasks and consumer applications, such as tag suggestion and organizing personal photo collections.

## References

[1] A. C. Berg, T. L. Berg, H. Daumé III, J. Dodge, A. Goyal, X. Han, A. Mensch, M. Mitchell, A. Sood, K. Stratos, and K. Yamaguchi. Understanding and predicting importance in images. In *Proc. CVPR*, 2012. 4

[2] O. Cappé and E. Moulines. On-line expectation-maximization algorithm for latent data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613, 2009. 4, 5

[3] K. Chatfield and A. Zisserman. Visor: Towards on-the-fly large-scale object category retrieval. In *Asian Conference on Computer Vision*, Lecture Notes in Computer Science. Springer, 2012. 2

[4] X. Chen, A. Shrivastava, and A. Gupta. Enriching visual knowledge bases via object discovery and segmentation. In *Proc. CVPR*, 2014. 2

[5] T.-S. Chua, J. Tang, R. Hong, H. Li, Z. Luo, and Y. Zheng. NUS-WIDE: a real-world web image database from National University of Singapore. In *Proc. CIVR*, 2009. 1, 2, 4

[6] A. Farhadi, I. Endres, D. Hoiem, and D. Forsyth. Describing objects by their attributes. In *Proc. CVPR*, 2009. 2

[7] R. Fergus, P. Perona, and A. Zisserman. A visual category filter for google images. In *Proc. ECCV*, 2004. 2

[8] Y. Gong, Y. Jia, T. K. Leung, A. Toshev, and S. Ioffe. Deep convolutional ranking for multilabel image annotation. In *Proc. ICLR*, 2014. 2, 6, 8

[9] Y. Gong, L. Wang, M. Hodosh, J. Hockenmaier, and S. Lazebnik. Improving image-sentence embeddings using large weakly annotated photo collections. In *Proc. ECCV*, 2014. 2

[10] J. Hoffman, E. Rodner, J. Donahue, K. Saenko, and T. Darrell. Efficient learning of domain-invariant image representations. In *Proc. ICLR*, 2013. 6

[11] S. J. Hwang and K. Grauman. Reading between the lines: Object localization using implicit cues from image tags. *PAMI*, June 2012. 2

[12] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding, 2014. http://arxiv.org/abs/1408.5093. 4, 6

[13] S. Karayev, M. Trentacoste, H. Han, A. Agarwala, T. Darrell, A. Hertzmann, and H. Winnemoeller. Recognizing image style. In *Proc. BMVC*, 2014. 2, 3

[14] A. Krizhevsky, I. Sutskever, and G. E. Hinton. ImageNet Classification with Deep Convolutional Neural Networks. In *NIPS*, 2012. 4

[15] L.-J. Li and L. Fei-Fei. OPTIMOL: Automatic object picture collection via incremental model learning. *IJCV*, 88:147–168, 2010. 2

[16] G. A. Miller. Wordnet: A lexical database for english. *Commun. ACM*, 38(11), 1995. 3

[17] N. Murray, D. Barcelona, L. Marchesotti, and F. Perronnin. AVA: A Large-Scale Database for Aesthetic Visual Analysis. In *CVPR*, 2012. 2, 3

[18] V. Ordonez, J. Deng, Y. Choi, A. C. Berg, and T. L. Berg. From large scale image categorization to entry-level categories. In *Proc. ICCV*, 2013. 4

[19] G. Patterson, C. Xu, H. Su, and J. Hays. The sun attribute database: Beyond categories for deeper scene understanding. *IJCV*, 108:59–81, 2014. 2, 3

[20] V. C. Raykar, S. Y, L. H. Zhao, G. H. Valadez, C. Florin, L. Bogoni, and L. Moy. Learning from crowds. *JMLR*, 11:1297–1322, 2010. 2, 4

[21] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge, 2014. http://arxiv.org/abs/1409.0575. 2, 3

[22] K. Saenko and T. Darrell. Unsupervised learning of visual sense models for polysemous words. In *Proc. NIPS*, 2008. 3

[23] W. J. Scheirer, N. Kumar, P. N. Belhumeur, and T. E. Boult. Multi-attribute spaces: Calibration for attribute fusion and similarity search. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2012. 6

[24] F. Schroff, A. Criminisi, and A. Zisserman. Harvesting Image Databases from the Web. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 33(4):754–766, apr 2011. 2

[25] P. Welinder, S. Branson, S. Belongie, and P. Perona. The Multidimensional Wisdom of Crowds. In *Proc. NIPS*, 2010. 2

[26] J. Xiao, J. Hays, K. Ehinger, A. Oliva, and A. Torralba. Sun database: Large-scale scene recognition from abbey to zoo. In *Proc. CVPR*, 2010. 2, 3

[27] B. Zhou, A. Lapedriza, J. Xiao, A. Torralba, and A. Oliva. Learning Deep Features for Scene Recognition using Places Database. In *Proc. NIPS*, 2014. 2, 3

[28] G. Zhu, S. Yan, and Y. Ma. Image tag refinement towards low-rank, content-tag prior and error sparsity. In *Proc. MM*, 2010. 2

## A. Details of correspondence calculation

Here we explain how we estimated the percentage of Flickr tags absent from ImageNet concepts (Section 3.2 of
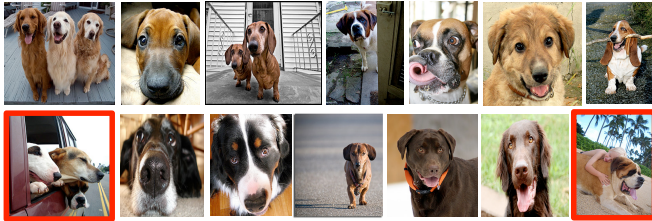
the submission). We collected the top 1000 Flickr tags, and manually filtered out non-image and location tags, with 612 tags remaining. We determined an automatic mapping from Flickr tags to WordNet synsets, by mapping each tag to its top WordNet noun synset, and manually corrected mismatches in the top 100 Flickr tags. We call an ImageNet synset *large* if it has 1000 of more node in the subtree. Of the top 100 Flickr tags, we found that 54 of them had large ImageNet subtree before correcting mismatches, and 62 had large subtrees after manual corrections. Of the remaining 512 uncorrected tags, 189 (37%) have large subtrees. Linear extrapolation suggests that $1 - ((189.0 * (62.0/54.0)) + 62)/612 = 54\%$ of tags are missing ImageNet subtrees. Of course, there are a number of questionable assumptions in this model, e.g., 1000 images may not be enough images for many classes, such as art.

## B. Issues with NUS-WIDE annotation

Figure 8 shows some samples of annotation error in the NUS-WIDE dataset. Another example is car and vehicle categories: in the NUS-WIDE test set there are 431 instances of "cars" of which only 177 instances are also annotated as "vehicle".



(a) Cat



(b) Dog

Figure 8: NUS-WIDE annotation error examples. The top retrieved images in RLR for cat and dog categories are shown. Red boxes are shown around images marked as negative samples in the dataset.

## C. Basic Logistic Regression

$$
\begin{array}{lcr}
\mathbf{w} & \text{logistic weights} & (14) \\
\mathbf{x} & \text{image features} & (15) \\
s = \mathbf{w}^T \mathbf{x} & \text{score given image data} & (16) \\
y \in \{0, 1\} & \text{observed label for each image} & (17)
\end{array}
$$

The model of label probabilities given image data is:

$$
\begin{aligned}
s &= \mathbf{w}^T \mathbf{x} & (18) \\
\sigma(s) &= \frac{1}{1 + e^{-s}} & (19) \\
P(y = 1|s) &= \sigma(s) & (20)
\end{aligned}
$$

The loss function for a dataset $\{(\mathbf{x}_i, y_i)\}$ is

$$
\begin{aligned}
L &= -\ln P(y_{1:N}|\mathbf{x}_{1:N}) & (21) \\
&= -\ln \left( \prod_{i:y_i=1} P(y_i = 1|\mathbf{x}_i) \right) \left( \prod_{i:y_i=0} P(y_i = 0|\mathbf{x}_i) \right) & (22) \\
&= \sum_i \left( -y_i \ln P(y_i = 1|\mathbf{x}_i) - (1 - y_i) \ln P(y_i = 0|\mathbf{x}_i) \right) & (23) \\
&= \sum_i \left( -y_i \ln \sigma(s_i) - (1 - y_i) \ln(1 - \sigma(s_i)) \right) & (24)
\end{aligned}
$$

We can also rearrange terms:

$$
\begin{aligned}
1 - \sigma(s) &= \frac{1 + e^{-s}}{1 + e^{-s}} - \frac{1}{1 + e^{-s}} = e^{-s}\sigma(s) & (25) \\
L &= \sum_i \left( -y_i \ln \sigma(s_i) - (1 - y_i) \ln e^{-s}\sigma(s_i) \right) & (26) \\
&= \sum_i \left( -\ln \sigma(s_i) + (1 - y_i)s_i \right) & (27) \\
&= \sum_i \left( \ln(1 + e^{-s}) + (1 - y_i)s_i \right) & (28)
\end{aligned}
$$

**Gradients.** During optimization, we use the gradients with respect to $\mathbf{w}$:

$$
\begin{aligned}
\frac{d}{d\mathbf{w}}\sigma(s) &= \sigma(s)\sigma(s)e^{-s}\mathbf{x} & (29) \\
&= \sigma(s)(1 - \sigma(s))\mathbf{x} & (30) \\
\frac{dL}{d\mathbf{w}} &= \sum_i \left( -\frac{y_i}{\sigma(s_i)}\frac{d}{d\mathbf{w}}\sigma(s_i) - \frac{1 - y_i}{1 - \sigma(s_i)}\frac{d}{d\mathbf{w}}(1 - \sigma(s_i)) \right) & (31) \\
&= \sum_i \left( -y_i(1 - \sigma(s))\mathbf{x}_i + (1 - y_i)\sigma(s_i)\mathbf{x}_i \right) & (32) \\
&= \sum_i \left( \sigma(s_i) - y_i \right) \mathbf{x}_i & (33)
\end{aligned}
$$

Note that this is zero when $y_i = \sigma(s_i)$, which indicates a perfect data fit.

Derivation using alternate form:

$$\frac{dL}{d\mathbf{w}} = \sum_i \left((\sigma(s_i) - 1)\mathbf{x}_i + (1 - y_i)\mathbf{x}_i\right) \quad (34)$$

$$= \sum_i (\sigma(s_i) - y_i)\mathbf{x}_i \quad (35)$$

## D. Robust Logistic Regression

$$\mathbf{w} \qquad \text{logistic weights} \quad (36)$$
$$\mathbf{x} \qquad \text{image features} \quad (37)$$
$$s = \mathbf{w}^T\mathbf{x} \qquad \text{score given image data} \quad (38)$$
$$y \in \{0, 1\} \quad \text{observed label for each image} \quad (39)$$
$$z \in \{0, 1\} \quad \text{hidden true label for each image} \quad (40)$$

The model of observations given scores is

$$P(z = 1|s) = \sigma(s) \quad (41)$$
$$P(y = 1|z = 1) = \pi \quad (42)$$
$$P(y = 0|z = 1) = 1 - \pi \ \text{false negative probability} \quad (43)$$
$$P(y = 0|z = 0) = \gamma \quad (44)$$
$$P(y = 1|z = 0) = 1 - \gamma \ \text{false positive probability} \quad (45)$$

In the paper, we fix $\gamma = 1$.

The marginal probability of a given observation is:

$$P(y|s) = \sum_{z \in \{0,1\}} P(y, z|s) = P(y|z = 1)P(z = 1|s) + P(y|z = 0)P(z = 0|s) \quad (46)$$

$$P(y = 1|s) = \pi\sigma(s) + (1 - \gamma)(1 - \sigma(s)) \quad (47)$$
$$= \pi\sigma(s) + (1 - \gamma)e^{-s}\sigma(s) \quad (48)$$
$$= \sigma(s)((1 - \gamma)e^{-s} + \pi) \quad (49)$$
$$P(y = 0|s) = (1 - \pi)\sigma(s) + \gamma(1 - \sigma(s)) \quad (50)$$
$$= (1 - \pi)\sigma(s) + \gamma e^{-s}\sigma(s) \quad (51)$$
$$= \sigma(s)(1 - \pi + \gamma e^{-s}) \quad (52)$$

The Maximum Likelihood loss function can be written:

$$L = \sum_i \left(-y_i \ln(\pi\sigma(s_i) + (1 - \gamma)(1 - \sigma(s_i)))\right.$$
$$\left. -(1 - y_i) \ln((1 - \pi)\sigma(s_i) + \gamma(1 - \sigma(s_i)))\right) \quad (53)$$
$$= \sum_i \left(-y_i \ln \sigma(s)((1 - \gamma)e^{-s_i} + \pi) - (1 - y_i) \ln \sigma(s_i)(1 - \pi + \gamma e^{-s_i})\right) \quad (55)$$
$$= \sum_i \left(-\ln \sigma(s_i) - y_i \ln((1 - \gamma)e^{-s_i} + \pi) - (1 - y_i) \ln(1 - \pi + \gamma e^{-s_i})\right) \quad (56)$$

When $s_i > \sim 35$, and thus $P(z = 1|s) \approx 1$, the summand should be implemented as:

$$-y_i \ln \pi - (1 - y_i) \ln(1 - \pi) \quad (57)$$

**Gradients.** During optimization, we could use the gradients with respect to $\mathbf{w}$:

$$\frac{dL}{d\mathbf{w}} = \sum_i \left(\sigma(s_i) - 1 - y_i \frac{-(1 - \gamma)e^{-s_i}}{(1 - \gamma)e^{-s_i} + \pi} - (1 - y_i)\frac{-\gamma e^{-s_i}}{1 - \pi + \gamma e^{-s}}\right)$$

$$= \sum_i \left(\sigma(s_i) - 1 - y_i \frac{-(1 - \gamma)}{(1 - \gamma) + \pi e^{s_i}} - (1 - y_i)\frac{-\gamma}{(1 - \pi)e^{s_i} + \gamma}\right)$$

(Dividing by $e^s$ is done for stability. The case where $s$ is very large should also be handled by a separate condition.)

We also wish to optimize with respect to the parameters $\pi$ and $\gamma$:

$$\frac{dL}{d\pi} = \sum_i \left(-y_i \frac{-e^{-s_i}}{(1 - \gamma)e^{-s_i} + \pi} - (1 - y_i)\frac{e^{-s_i}}{1 - \pi + \gamma e^{-s_i}}\right) \quad (60)$$

$$= \sum_i \left(-y_i \frac{-1}{1 - \gamma + \pi e^{s_i}} - (1 - y_i)\frac{1}{(1 - \pi)e^{s_i} + \gamma}\right) \quad (61)$$

$$\frac{dL}{d\gamma} = \sum_i \left(-y_i \frac{-e^{-s_i}}{(1 - \gamma)e^{-s_i} + \pi} - (1 - y_i)\frac{e^{-s_i}}{1 - \pi + \gamma e^{-s_i}}\right) \quad (62)$$

$$= \sum_i \left(-y_i \frac{-1}{1 - \gamma + \pi e^{s_i}} - (1 - y_i)\frac{1}{(1 - \pi)e^{s_i} + \gamma}\right) \quad (63)$$

### D.1. Stochastic EM algorithm

In the E-step, we compute the probabilities over the latent $z$'s given the data and the current model.

$$\alpha_i \stackrel{\triangle}{=} \bar{P}(z = 1|y_i, s_i) = \frac{P(y_i|z = 1, s_i)P(z_i = 1|s_i)}{P(y_i|s_i)} \quad (64)$$

which is computed with

$$P(z = 1|y_i = 1, s_i) = \frac{\pi\sigma(s_i)}{\pi\sigma(s_i) + (1 - \gamma)(1 - \sigma(s_i))} \quad (65)$$

$$P(z = 1|y_i = 0, s_i) = \frac{(1 - \pi)\sigma(s_i)}{(1 - \pi)\sigma(s_i) + \gamma(1 - \sigma(s_i))} \quad (66)$$

**M-step derivation.** In the M-step, we update the various model parameters. It can be derived by minimizing the negative expected complete log-likelihood:

$$E = <-\sum_i \ln P(y_i, z_i|s_i)>_{\alpha_i} \quad (67)$$

$$= <-\sum_i \ln P(y_i|z_i)P(z_i|s)>_{\alpha_i} \quad (68)$$

$$<\ln P(y_i = 1|z_i)> = \alpha_i \ln \pi + (1 - \alpha_i) \ln(1 - \gamma) \quad (69)$$
$$<\ln P(y_i = 0|z_i)> = \alpha_i \ln(1 - \pi) + (1 - \alpha_i) \ln \gamma \quad (70)$$
$$<\ln P(z_i|s_i)> = \alpha_i \ln \sigma(s_i) + (1 - \alpha_i) \ln(1 - \sigma(s_i)) \quad (71)$$

$$(72)$$

The derivatives are then:

$$\frac{dE}{d\pi} = -\sum_i \left( y_i \frac{\alpha_i}{\pi} + (1 - y_i) \frac{-\alpha_i}{1 - \pi} \right) \quad (73)$$

$$\frac{dE}{d\gamma} = -\sum_i \left( -y_i \frac{1 - \alpha_i}{1 - \gamma} + (1 - y_i) \frac{1 - \alpha_i}{\gamma} \right) \quad (74)$$

$$\frac{dE}{d\mathbf{w}} = -\sum_i \left( \alpha_i (1 - \sigma(s_i)) - (1 - \alpha_i)\sigma(s_i) \right) \mathbf{x}_i \quad (75)$$

$$= \sum_i \left( \sigma(s_i) - \alpha_i \right) \mathbf{x}_i \quad (76)$$

Solving for $dE/d\pi = 0$ and $dE/d\gamma = 0$ gives:

$$\pi \leftarrow \frac{\sum_i y_i \alpha_i}{\sum_i \alpha_i} \quad (77)$$

$$\gamma \leftarrow \frac{\sum_i (1 - y_i)(1 - \alpha_i)}{\sum_i (1 - \alpha_i)} \quad (78)$$

**Stochastic EM algorithm.** In the stochastic EM algorithm, we keep running tallies of

$$S_y = \sum_i y_i / N \quad (79)$$

$$S_{y\alpha} = \sum_i y_i \alpha_i / N \quad (80)$$

$$S_\alpha = \sum_i \alpha_i / N \quad (81)$$

and then, in each M-step update, update the parameters as:

$$\pi \leftarrow \frac{S_{y\alpha}}{S_\alpha} \quad (82)$$

$$\gamma \leftarrow \frac{1 - S_y - S_\alpha + S_{y\alpha}}{1 - S_\alpha} \quad (83)$$

We can also use $dE/d\mathbf{w}$ as a gradient estimate instead of $dL/d\mathbf{w}$.