
A Cascade Ranking Model for Efficient Text to Image Retrieval

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 The ubiquitous accessibility and explosive increase of image data has resulted in
2 the property of research in image retrieval. In this paper, we address the task of
3 text-based image retrieval, to select a set of images from a large image database.
4 We vectorized the query descriptions to a sparse binary matrix, employed PLSRe-
5 gression to learn the common subspace between the visual and textual features,
6 made distance-based relevant image predictions, and rank the relevant images
7 based on a combination of predicted l2-distance and category jaccard similarity.
8 Experimental results demonstrate that our method effectively utilizes both visual
9 and semantic meanings of the cross-modal information, outperforming other teams
10 on Kaggle competition, and can exploit large scale vision and language datasets
11 for knowledge transfer.

12 1 Introduction

13 The task of image retrieval means that, if given a specific description, retrieving all images relevant
14 to that description within a potentially very large database of images. As multimedia data such as
15 texts and images are growing exponentially on digital devices, it becomes increasingly difficult for
16 users to search information effectively and efficiently. Therefore, cross-modal retrieval has attracted
17 a considerable attention from many researchers. The current research effort is to design variant
18 approaches to solve the cross-modal challenges more accurate and more scalable.

19 1.1 Problem Definition

20 As part of the Cornell Tech machine learning community, we are assigned with the task to build a
21 search engine using novel statistical learning techniques with a goal to retrieve images that share the
22 same semantics with the given query description.

23 **Data** During training, we have a dataset of 10,000 image samples, each with the following data
24 available for learning: a 224x224 JPG image, a list of tags indicating objects appeared in the image, a
25 five-sentence description and feature vectors extracted using ResNet, a state-of-the-art Deep-learned
26 CNN. The features are extracted from pool5 and fc1000 layer.

27 During testing, your system matches a single five-sentence description against a pool of 2,000
28 candidate samples from the test set. Each sample has: a 224x224 JPEG image, a list of tags for that
29 image, ResNet feature vectors for that image.

30 **Output** For each description, our system ranks each testing image with the likelihood of that image
31 matches the given sentence. Our system then returns the name of the top 20 relevant images.

32 **Evaluation Metric** We need to rank-order test images for each test description and the rating system
33 uses MAP@20 as the evaluation metric. If the corresponding image of a description is among our

34 algorithm's 20 highest scoring images, this metric would then give a certain score based on the
35 ranking of the corresponding image.

36 **2 Model Overview**

37 In this section, we describe our model for text to image retrieval and the training procedure in details.
38 At test time, 2,000 images and 2,000 natural language textual descriptions are provided. For each
39 description, the system needs to select from the 2,000 test images a subset of 20 target images that
40 match the query text and sorts the candidates by relevancy ranking.

41 Figure 1 shows the architect of our final learning model. Before reaching the final model, we tried
42 different approaches to analyze the cross-modal relationship between the image and the text. For
43 each approach, we tried to utilize different features and different parameters for fine-tuning.

44 We divided the problem into the follwing steps: 1). extracting features from training descriptions,
45 2). finding correlations between text features and image features, 3). combining the components
46 above to produce l2-distances for all 2,000 testing images, 4) re-ranking the images by combining
47 l2-distance and jaccard similarity then retrieving the top 20 candidates for a given text query.

48 **2.1 Extracting feaatues from descriptions**

49 **2.1.1 Attempt One: Doc2Vec and Word2Vec**

50 Doc2Vec is an unsupervised algorithm that generate vector for paragprahs in documents. The feature
51 vectors generated by Doc2Vec can be utilized to find similarity between documents. Upon given
52 the data, we quickly realized that we need to identify similar descriptions to perform the task given.
53 Yet it performs poorly. Our observation shows that although Doc2Vec is able to classify similar
54 documents, it cannot distinguish different nouns that captrues important meaning in the descriptions.
55 For example, Doc2Vec would give a similar score to "A man walks down the street" and "A woman
56 walks down the street". Therefore, it is not ideal for the given task.

57 We then turned to Word2Vec, which takes an average of the sum of the word vectors to generate
58 a sentence vector. The performance is not that ideal either. We suspect the reason is that the
59 training sample is too small. For example, Google GloVe is an unsupervised learning algorithm that
60 performs exceptionally for obtaining vector representations of words. Yet the training is performed
61 on aggregated global word-word co-occurence statistics.

62 **2.1.2 Attempt Two: TF-IDF**

63 TF-IDF refers to "Term Frequency-Inverse Document Frequency", and the TF-IDF weight is used to
64 evaluate how important a word is in a document. A term frequency(TF) is a count of how many times
65 a word occurs in a given document(like bag of words), and the inverse document frequency(IDF) is
66 the number of times a word occurs in a corpus of documents. Words that are used frequently in many
67 documents will have a lower weighting while infrequent ones will have a higher weighting. The final
68 weight is the product of TF and IDF, by calculating the final weight of given words in documents, we
69 can form the vector for every description and then use the vector to compute the similarity.

70 Firstly, we notice that every description is in natural language. Therefore, we pre-process the
71 descriptions by removing all punctuation, using NLTK module to remove all stopwords in the text,
72 convert all letters in the descriptions into lower letter. After preprocessing, we processed descriptions
73 corresponding to a certain image into Bag of Words, and we calculate the value vector for every
74 description file. We then calculate the TF-IDF array for each tag file after doing similar preprocessing.
75 By using vectors of descriptions and tags, we calculate cosine similarity between vectors. We ran a
76 naive experiment to predict testing results and used it as a baseline to check our performance.

77 **2.1.3 Attempt Three: SVM Tag Generation**

78 We did some exploratory data analysis so that we found that there are 80 different tags for training
79 set. Therefore, we trained 80 One vs. All SVM classifier with hinge loss and tried to predict if the
80 predicted tags are in a given description. Experimental results shows that the following text feature

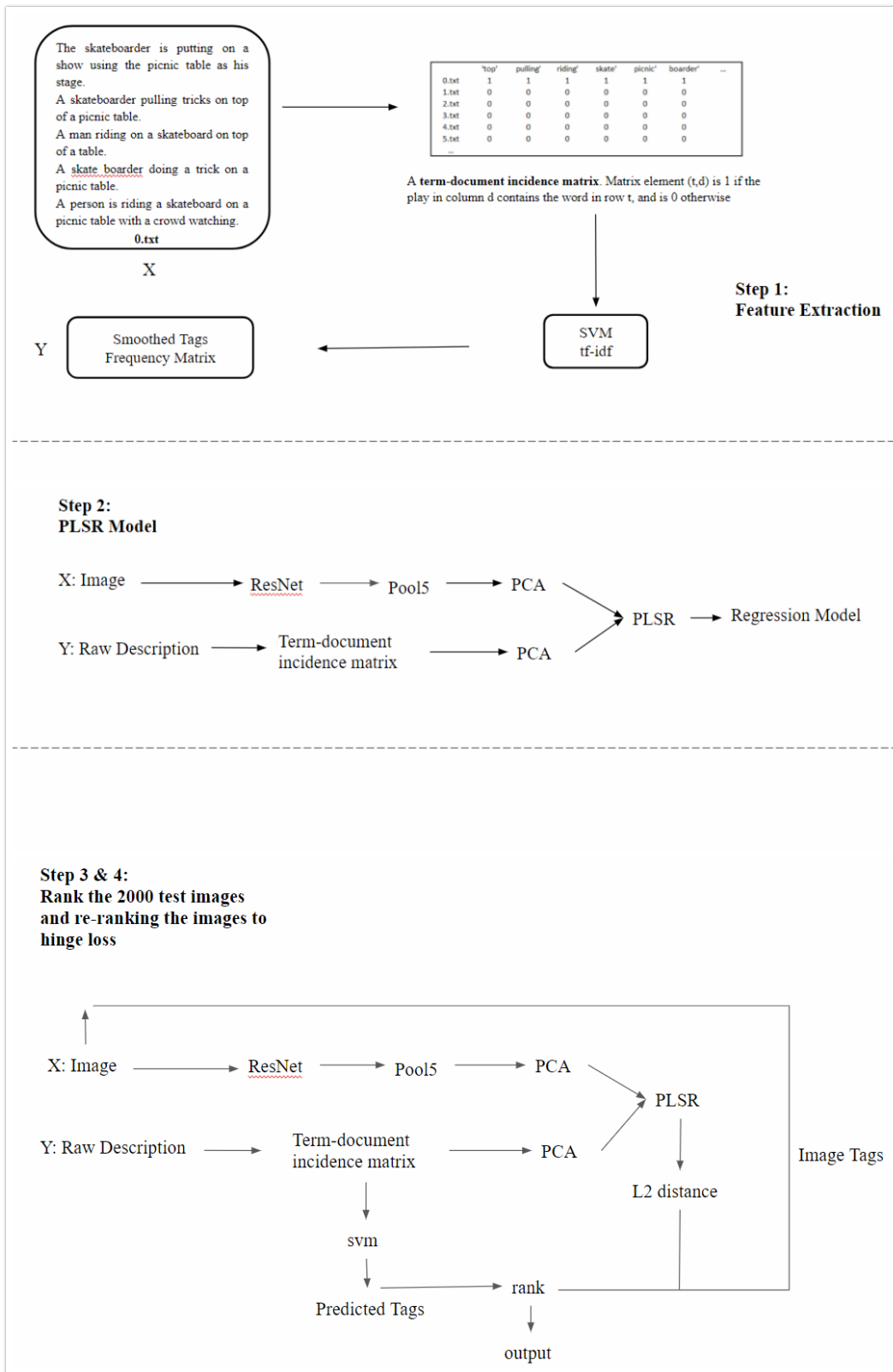


Figure 1: Model Architecture

81 extraction performs better than our SVM tag generation. Yet we were able to utilize these models
82 generated to obtain results that subsequently serves as input for our re-ranking algorithm in step 4.

83 2.1.4 Final Version: Term-document incidence matrix

84 Our final attempt is the simplest form of feature extraction for the text description. That is to create a
85 term-document incidence matrix. First, we scan through all the documents to find a set of unique
86 words after removing the stopwords that appears in the corpus. The set of unique words then becomes
87 a feature vector for each query description. The matrix element is 1 if a given word appears in the
88 description, and is 0 otherwise.

89 Instead of generating a sparse binary matrix, we have also tried recording the counts and thereby a
90 document term matrix. Experimental result shows the incidence matrix works better in our case as the
91 document term matrix would give an unfair weight to two documents if the overlapping term appears
92 multiple times when comparing similarity. In reality, if two descriptions have many intersections
93 yet the intersections are pointing to a single word, it is a weaker signal of similarity than if two
94 descriptions have less intersections yet each intersection represents an unique word. We have also
95 tried recording the counts of unique nouns that appears in the corpus. Yet experimental result shows
96 that incidence matrix with all the words except stop words works the best with our model.

97 We then utilized the term document incidence matrix as our textual features to map the relationship
98 with the visual features. Detail of the model is given below. For simplicity, we will refer to term
99 document incidence matrix(TDIM) as text feature below.

100 2.2 Finding correlations between image features and text features

101 We have tried to use textual features alone, including text features from descriptions and tags features
102 generated by our SVM model, to make predictions. It is foreseeable that the accuracy would not be
103 ideal. Since the textual features are manually annotated, descriptions and tags usually do not consist a
104 comprehensive information of the image. As a lot of information is lost in this way, it is necessary to
105 incorporate the image features (fc1000 and pool5 from the ResNet) and construct a mapping between
106 text features and images features.

107 2.2.1 Step One: Principal Component Decomposition

108 In large scale data mining and predictive modeling, especially for multivariate regression, we often
109 start with a large number of possible explanatory/predictive variables. Therefore, variable selection
110 and dimension reduction is a major task before modeling our cross-modal relationship.

As our text feature is represented by a sparse binary matrix, it makes sense to use principal component analysis to perform linear dimensionality reduction using Singular Value Decomposition of the data to project it to a lower dimensional space. PCA is also a powerful tool for pattern recognition which is often used in image processing. Therefore, we perform reduction technique on both features before training our model that maps one to the other.

$$\begin{aligned} pool5 &\rightarrow PCA(n - components = 1000) \rightarrow Reduced_pool5 \\ TDIM &\rightarrow PCA(n - components = 3500) \rightarrow Reduced_TDIM \end{aligned}$$

111 2.2.2 Step Two: Partial Least Square Regression

112 After some literature review, we found some popular methods for establishing inter-modal relationships
113 between data from different modalities, including Canonical Correlation Analysis, Bilinear Model, and
114 Partial Least Squares. After many experimentations, we decided to use PLSRegression that linearly
115 map images with different modalities to a common subspace. According to Chen et al., they apply
116 PLS to image features into the text spaces in multimedia retrieval, then learn a semantic space for
117 the measure of covariance between the two modalities. In PLSRegression, it maximizes correlation
118 between the input and output as well as maximizes fit while minimizing misfits.

119 We used different combination of X and Y to find a model that best explains the information presented
120 by the given dataset. For Y value, we used all the textual features mentioned in Section 2.1. For X
121 value, we conducted different trials with fc1000, pool5, and stacking both fc1000 and pool 5 as our

input feature. After experimenting with different combinations, we identified that term document incidence matrix and pool5 features gives us the best accuracy with the model proposed.

$$pool5 \rightarrow PLSR(n - components = 800) \rightarrow Reduced_TDIM$$

2.3 Calculating l2-distances for all 2000 images

After we trained the PLSRegression model with features and parameters listed above, we then used the end result of regression to calculate euclidean distance between each predicted result and the true description features across the 2,000 test images. Then we rank the 2,000 test images based on the l2 distance in an ascending order to obtain a preliminary ranking results for our test set descriptions.

2.4 Re-ranking the images

We then introduce another similarity measure to help us capture the categorical connection between query descriptions and images. The similarity is defined as follows:

$$sim = \frac{|D| \cap |I|}{|I|}$$

Where D is the set of predicted tags for query document and I is the set of tags for image candidates. Finally, for a give query document D we will reward pair-wise l2 distance from 2.3 and pick the top 20 as our output:

$$min_{l2}\{l2_{ij} - sim_{ij}\}$$

3 Experiments

We conducted multiple experiments to blend our features and models in a way that boosts our performace. After a number of trials, we obseved that with the same model, using pool5 features as our X for PLSRegression performs much better than fc1000. This makes sense because the 1000-dimensional layer is the one just before the softmax, so the i dimesion captures how confident the network is that the image contains a class i. In other words, fc1000 provides categorical information of the image. In comparision, the pool5 layer captures some kind of semantic meaning, although each dimension is not necessarily interpretable. In comparision, pool5 has more images features and may share the same semantics meaning with the text query. We suspects that might be why it outperforms fc1000 by a significant amount in our model setup.

As for our text features, we find the less processing, the better it works in terms of letting image features switch into its space in the subsequent regression analysis. The resulted common semantic space provides a good measure of correlation between the two different modalities.

3.1 Experimental Results

Figure 2 displays a list of testing accuracy from our selected testing trials. Our best result achieved on the private leader board is 0.42983.

4 Conclusion

In this work, we demonstrated a two stage text-based image retrieval with a re-ranking method that significantly boosts the performance of our retrieval task (see Figure 3). Cross-modal retrieval is a more powerful and effective way to multimodal retrieval than the traditional single modality based techniques. In the proposed model, we leveraged the text annotations to learn a visual embedding of the images and showed that this embedding predicts very well when given a text query (see Figure 4 and 5). In addition to locate the corresponding image, our model is also able to group images that share the same semantic meanings. Finally, a bottleneck of our performance is due to the inconsistency in manual annotations, as there are 84 training images that lack of tag descriptions. Therefore, our joint textual and visual model can leverage refining dataset to provide a more accurate method to query large image database.

Model Used	Accuracy on Test Set
TF – IDF	0.14795
SVM tags	0.15024
Word2Vec	0.16075
$PLSR_{W2V_fc1000}$	0.16430
$PLSR_{TDIM_{NOUNS} \rightarrow fc1000}$	0.25095
$PLSR_{TDIM_{NOUNS} \rightarrow pool5}$	0.32139
$PLSR_{TDIM \rightarrow pool5}$	0.34639
$PLSR_{pca+pool5+svm+naive(euc-jaccard/min_dist)}$	0.35441
$PLSR_{pool5+svm+naive(euc-jaccard(len(test_img))/min_dist)}$	0.38831
$PLSR_{pool5+svm+naive(euc-jaccard(len(test_img))/min_dist)}$	0.42983

**The disparity between the last two entries is due to the difference of n components。
 **The lower one is obtained by using n components = 800
 **The higher score is obtained by using n components = 1000

Figure 2: A selection of our submitted Kaggle results



Figure 3: Sample Re-ranking Performance

4.1 Future Works

Due to the limited amount of time and computation resources we have, we can only run a small number of experiments before the deadline. We believe that we can improve our model by leveraging the fc1000 image features and combining it with the existing model to build an ensemble model. Lastly, non-linear modeling might be an area worth exploring, as literature review shows that methods such as Deep Canonical Correlation Analysis(DCCA) and Multimodal Latent Binary Embedding(MLBE) seems promising.

Acknowledgments

The authors would like to thank Serge Belongie for the fantastic lectures. The authors would also like to thank Michael Wilber and Xun Huang for the insightful replies on Slack.

A pack of zebra walking along a grassy field.
 several wild zebra's standing in a field of
 grasses and shrubs
 A bunch of zebras standing in a field
 A pack of zebras stand on the plain
 looking into the distance
 A herd of zebras are standing on a
 grassy field.



Text Query

Retrieved Images

Figure 4: Sample Query Result

A man in a baseball uniform
 kneeling down while holding a
 baseball bat.
 This ball player is posing for a
 picture in his Ranger's uniform.
 A baseball player wears a Rangers
 outfit while kneeling and holding a
 bat.
 A baseball player poses for a picture
 with his baseball bat.
 Baseball player posing in his
 uniform holding a baseball bat.



Text Query

Retrieved Images

Figure 5: Sample Query Result

167 **References**

- 168 [1] Albert Gordo & Diane Larlus (2017) *Beyond instance-level image retrieval: Leveraging captions to learn a*
169 *global visual representation for semantic retrieval*
- 170 [2] Christopher D. Manning & Prabhakar Raghavan & Hinrich Schütze (2008) An example information retrieval
171 problem *Introduction to Information Retrieval*
- 172 [3] Wengang Zhou & Houqiang Li & Qi Tian (2017) *Recent Advance in Content-based Image Retrieval: A*
173 *Literature Survey*
- 174 [4] Shangwen Li & Sanjay Purushotham & Chen Chen & Yuhuo Ren (2017) *Measuring and Predicting Tag*
175 *Importance for Image Retrieval*
- 176 [5] Jeffrey Pennington & Richard Socher & Christopher D. Manning *GloVe: Global Vectors for Word Represen-*
177 *tation*
- 178 <https://docs.scipy.org/doc/numpy/>
- 179 <http://scikit-learn.org/>
- 180 <https://docs.python.org/3/library/pickle.html>
- 181 <https://pypi.python.org/pypi/nltk>