

Recommending Outfits from Personal Closet

Pongsate Tangseng¹, Kota Yamaguchi², and Takayuki Okatani^{1,3}

¹Tohoku University

²CyberAgent, Inc.

³RIKEN Center for AIP

¹{tang seng, okatani}@vision.is.tohoku.ac.jp

²yamaguchi_kota@cyberagent.co.jp

Abstract

We consider grading a fashion outfit for recommendation, where we assume that users have a closet of items and we aim at producing a score for an arbitrary combination of items in the closet. The challenge in outfit grading is that the input to the system is a bag of item pictures that are unordered and vary in size. We build a deep neural network-based system that can take variable-length items and predict a score. We collect a large number of outfits from a popular fashion sharing website, Polyvore, and evaluate the performance of our grading system. We compare our model with a random-choice baseline, both on the traditional classification evaluation and on people’s judgment using a crowdsourcing platform. With over 84% in classification accuracy and 91% matching ratio to human annotators, our model can reliably grade the quality of an outfit. We also build an outfit recommender on top of our grader to demonstrate the practical application of our model for a personal closet assistant.

1. Introduction

There have been growing interests in applying computer vision to fashion, perhaps due to the rapid advancement in computer vision research [1–10]. One of the popular fashion applications is item recommendation [11–14], where the objective is to suggest items to users based on user’s and/or society’s preference. Computer vision is used in various fashion applications such as e-commerce and social media. Recently, Amazon announced their automatic style assistant called “Echo Look™”. Although the underlying mechanism is not published, emerging commercial applications confirm the ripe of computer vision applications in fashion.

Measuring the quality of outfit is essential in building



Figure 1: Given an arbitrary number of items, our goal is to evaluate the quality of the outfit combination.

fashion recommendation system. In this paper, we consider the problem of grading arbitrary combination of items as a whole (Figure 1). Previous works in outfit evaluation can be divided into two groups based on the input format: a worn outfit as a full-body picture as in [3, 15–17], and as a set of images of items [12, 14], or a combination of both [11]. Each outfit can have an arbitrary number of items. For examples, in one day, one might prefer a combination of a jacket, a t-shirt, and jeans, while in the another she might want to wear a dress. Our goal is to build a machine learning system that accepts a variable numbers of items yet produce a consistent score for any size of combinations.

In this paper, we view an outfit as a bag of fashion items and utilize deep neural networks to produce a score for a fixed-length representation of outfits. Unlike style recognition [3, 10, 15, 17], we take item images in isolation, not on human body, as seen on e-commerce sites or catalogs. We collect a large number of outfit data from a popular fashion website polyvore.com, and evaluate our approach based on standard classification metrics and human judgment. Our Polyvore409k dataset consists of 409,776 sets of clothing items from 644,192 unique items. The dataset forms a large bipartite graph of items and outfits. We partition the dataset into training and testing sets such that there is no overlapping nodes and edges between the sets, and use them measure the classification performance. We also conduct a human study using crowdsourcing to assess predicted scores against human judgment, and show our model closely resembles human behavior. Using our grader, we

build an outfit recommendation system that takes clothing items as an input and recommends the best outfits from the given items, to demonstrate the usefulness in a real-world scenario of personal outfit assistant. The contributions of the paper are summarized below:

1. We build Polyvore409k dataset containing 409,776 outfits and 644,192 items. Every outfit covers the entire body with a variable numbers of items.
2. We propose an outfit grader that produces a score for fashion outfits with a variable number of items. Our empirical study shows that our model achieves 84% of accuracy and precision in Polyvore409k dataset.
3. We propose a human judgment framework on outfit quality, which provides a simple and reliable method to verify the reliability of outfit verifiers using a crowdsourcing platform.
4. We demonstrate that our outfit grader can build a recommendation system that suggests good outfits from a pool of items.

2. Related Work

2.1. Outfit Modeling

The use of computer vision techniques to study fashion is gaining popularity. Some early studies work on outfit images [15–17]. Although these studies can use the appearances of outfit on a real subject, accurately identifying items in an outfit image is still an open problem. The manual annotation is costly, and the automatic detection and segmentation of fashion items in an outfit image [1, 2, 4–6, 18] are still not reliable. For example, “dress” and “top with skirt” are often incorrectly segmented, and the small objects like shoes and accessories are often missed. In addition, the importance of each item to the overall style may or may not be related to the scale of the item in an outfit image. In this paper, we aim to study outfits as a combination of items, where each item has its own image.

Some studies treat outfits as combinations of item images as well [11, 12, 14, 19]. In [11], items in an outfit are recommended according to the requested occasion and the existing items in that outfit. The work by [12] focuses on learning personal preference on fashion based on accounts and associated outfits from polyvore.com. A study of pairwise relationship between fashion items was explored in [19] using co-purchase data. The work by [14] also uses data from polyvore.com to learn outfits as combinations of items based on item image, name, and category. They create an item recommendation system that suggests an item to match with other manually selected items.

Outfits have a natural structure based on human body, but [11, 12, 14] consider outfits with fixed number of items without considering variation in the structure. Outfits in [12] consist of one top, one bottom, and a pair of shoes without

considering full-body items such as a dress, nor accessories. In [11], recommendation is made for either whole outfit, or upper-lower body pairs. Likewise, outfits in [14] consist of 4 items, regardless of item role. [14] does not guarantee the completeness of outfits. Since an outfit can be a collection of any items, it is possible to have incomplete outfits that do not cover whole body, such as an outfit consisting only of four pairs of boots.

In this work, we view outfits as collections of items from polyvore.com, similar to [12, 14]. We arrange the outfit data such that each outfit covers the entire body by considering the body part covered by each item, with variable number of items in the outfit.

2.2. Fashion Datasets

The number of fashion datasets is growing. Each image in some datasets [2, 4, 6, 16] is an outfit images, which contains many items. In other datasets [12–14, 19], each image contains only one item. Some datasets [11, 20] are combinations of both type. In [19], item combinations come in pairwise format from amazon co-purchase data. However, items that are bought together do not necessarily mean that they look good together as an outfit.

There are segmentation datasets [1, 4, 6] that seem suitable for our problem setup, because the datasets provide outfit images with the boundary of each item, but the number of samples is too few to learn a reasonable predictive model. Although [12] and [14] use datasets with combination of images as outfits and each item has its own image, the dataset is not publicly available. For the above reasons, we collect and build a new dataset, Polyvore409k dataset, which we describe in section 3.

In fashion outfit problem, each sample is an outfit which is a combination of items. We have to cleanly separate training data from testing data both for a set and individual items. [11] and [12] do not describe the detail on separation. [14] constructs a graph dataset, where each node represents an outfit, and a connection between any two nodes is formed if these two outfits have a common items. After that, the graph is segmented based on connected components. In this work, we use an efficient alternative approach to split a graph, which we describe in section 3.4.

3. Polyvore409k Dataset

This section describes our Polyvore409k dataset that consists of variable-length sets of items. Our Polyvore409k dataset has 409k outfits consisting of 644k item images. The comparison of outfit datasets to the previous work is shown in table 1. We plan to release the metadata of items and outfits, including the URLs to the images, to the public.

Table 1: Comparison of outfit datasets

	[11]	[12]	[14]	Ours
#images in dataset	24,417	85,252	347,339	644,192
Annotation method	Crowdsourcing	Metadata from polyvore.com		
Outfit labels	Occasions and attributes	User-created: positive Randomly created: negative	Based on votes	User-created: positive Randomly created: negative
#items in outfits	2	3	4	Variable, up to 8
Human body parts	2 parts	3 parts	No	6 parts
Train/test item separation	Not verified		Verified	
Evaluated by human	No	No	No	Yes

Table 2: Number of unique items in each outfit part

Part	Outer	Upper	Lower	Full	Feet	Accessory
Train	11,168	21,760	16,287	11,523	26,574	60,760
Test	6,656	12,744	11,089	8,871	17,564	37,988

3.1. Data Collection

We collect Polyvore409k dataset from the fashion-based social media website polyvore.com. Each outfit, or *set* in Polyvore’s terminology, consists of a title, items in the set, a composed image, and behavioral data such as likes and comments from other users.

3.2. Data Preprocessing

Data Cleansing The collected sets can contain non-clothing items or items that cannot be worn such as logo, background image for presentation purpose, or cosmetic items. We remove the item if its name does not contain clothing categories. After that, each item in a set is categorized into one of 6 outfit parts according to its categories:

1. Outer: coat, jacket, parka, etc.
2. Upper-body: blouse, shirt, polo, etc.
3. Lower-body: pants, jeans, skirt, joggers, etc.
4. Full-body: dress, gown, jumpsuit, robe, etc.
5. Feet: shoes, boots, flats, clutches, etc.
6. Accessory: bag, glove, necklace, earring, etc.

Our definition of an outfit is a set that covers both upper and lower part of body, each of first 5 categories has at most one item, and at most three items for accessory. Sets that do not cover the whole body, e.g. missing lower body, are removed. At the end, we obtained 409,776 valid outfits which are composed of 644,192 unique items.

We consider only two layers on the upper body (outer and upper) because of the visibility, as the layers under two outermost layers are usually covered. We process sweaters, knitwear, and the likes as outer-upper hybrids. They will be considered as an upper if the outfit has other outer, and as an outer if the outfit does not have. The list of item categories and respective outfit parts is included as supplementary.

3.3. Quality Measurement

Measuring the quality of an outfit is a challenging task due to the subjective nature of judging visual appearance. The approach of [14] directly uses the number of votes (or *like* in polyvore.com’s terminology) of the outfit on the website as a quality measurement. However, some studies [16, 17, 21] argue that the number of votes from social media does not directly reflect the quality of the outfit, because the number of clicks is affected by a variety of factors, such as the topology of the social networks or the time when the outfit was published. In [12], the quality is defined by the preference of each user: outfits created by the user are treated as positive samples, and outfits created by randomly pick items are treated as negative. Given these insights, we take the following strategy.

Positive Samples Each Polyvore409k outfit has an associated *likes* that Polyvore users provide. Although the number of *like* might not directly reflect the quality of the outfit, it still shows that some people like the outfit. As a result, we use 212,623 outfits that has least one *like* as positive samples. In the future, as we obtain more data, we wish to increase the number of *like* threshold.

Negative Samples Similar to [12], we use outfits created by picking items randomly as negative samples. We believe that there are some preferred combinations of colors, textures, or shape of items in an outfit, and we assume that a randomly created outfit has very small chance to match those preferences.

For each positive sample, we create two identical samples as negative samples, because the number of preferred combinations is expected to be significantly lower than random combinations. Then, we replace items in those two negative samples with random items of the same parts from the same train/test item pool. Although this sampling strategy is not *i.i.d.*, this approach guarantees the disjoint set property between training and testing sets, and tends to produce *hard* negative examples that shares some items with positive counterpart. Also, the distribution of number of items and existences of outfit parts in samples are preserved.

Algorithm 1: Disjoint Set Sampling

```

input : All outfits  $O$ 
output: Set  $A$ ,  $B$ , and  $C$  containing outfits such that items in outfits
         in  $A$  is not in  $B$  and vice versa
 $A \leftarrow B \leftarrow C \leftarrow \emptyset$ ;
 $A \cup \{O_0\}$ ;
for  $i \leftarrow 1$  to  $|O|$  do
     $O \leftarrow O_i$ ;
     $items_A \leftarrow$  items in outfits in  $A$ ;
     $items_B \leftarrow$  items in outfits in  $B$ ;
     $items_O \leftarrow$  items in  $O$ ;
     $sec_{AO} \leftarrow intersection(items_A, items_O)$ ;
     $sec_{BO} \leftarrow intersection(items_B, items_O)$ ;
    if  $|sec_{AO}| > 0$  and  $|sec_{BO}| > 0$  then  $C \cup \{O\}$ ;
    else if  $|sec_{AO}| > 0$  then  $A \cup \{O\}$ ;
    else if  $|sec_{BO}| > 0$  then  $B \cup \{O\}$ ;
    else
        if  $|A|/2 > |B|$  then  $B \cup \{O\}$ ;
        else  $A \cup \{O\}$ ;
    end
end

```

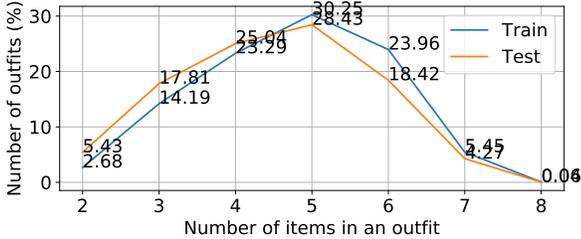


Figure 2: Distribution of number of items in outfits in training and testing partition are similar.

Table 2 shows the number of items in each part of outfit. Figure 2 shows the distribution of number of items in an outfit in train and test splits for both positive and negative samples. Table 3 shows the numbers of positive and negative samples in each split.

3.4. Evaluation Data

The set-item relationship constitutes a bipartite graph, where nodes are outfits or items, and edges represent inclusion relationship. For performance evaluation using Polyvore409k, we have to split the bipartite graph such that the training and testing splits do not share any item or outfits. We use Algorithm 1 to separate training and testing splits.

4. Outfit Grader

4.1. Problem Formulation

We formulate the outfit grading as a binary classification problem. Given an outfit $O \equiv \{x_{outer}, x_{upper}, \dots, x_{accessory3}\}$, where x_{part} is an item image, the goal is to learn a mapping function: $F : O \mapsto y$ to predict the outfit’s quality $y \in \{0, 1\}$.

Table 3: Number of outfits in each train and test partition

Number of outfits	Train	Test
Positive samples	66,434	26,813
Negative samples	132,868	53,626
Total	199,302	80,439
Ratio positive:negative	1:2	1:2

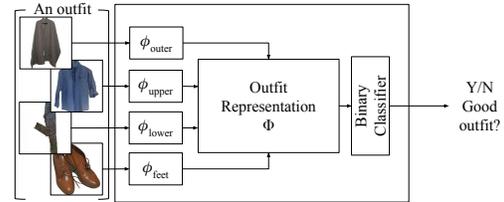


Figure 3: Outfit Grader

Once we learn the mapping F , we are able to sort arbitrary combinations of items according to the prediction score.

The challenge is how to represent an outfit O with a variable number of items. Luckily, the number of visible items is limited even though an outfit can contain a variable number of items. Therefore, we assign each item into one of the six categories and concatenate the item representations to produce the outfit representation. Figure 3 shows our grader. Our grader takes a bag of images and convert them to feature representations, then concatenates the individual features according to the item’s category to produce the fixed-length representation. We describe details below.

4.2. Item Representation

We convert the image of each item in the outfit to a feature representation $\phi_{part}(x_{part})$, using a convolutional network. In this paper, we use ImageNet-pretrained ResNet-50 [22], and extract the 2,048-dimensional embedding from `pool5` layer as an item representation. We extract features for 5 item parts and up to 3 accessories. For missing parts, we give a mean image to obtain features which is equal to zero-input to the convolutional network.

4.3. Outfit Representation

After we extract features from each item, we concatenate all features in the fixed order to form an outfit representation $\Phi(O) \equiv [\phi_{outer}, \phi_{upper}, \dots, \phi_{accessory2}]$. Note that we allow accessories to appear multiple times in the outfit, and we simply concatenate all the accessory features ignoring the order. Outfits with less than 3 accessories get mean images as well to the other part. We have 5 item parts and 3 accessories per outfit, resulting in a 16,384 dimensional representation as the outfit representation.

Table 4: Accuracy, average precisions, and average recall of our outfit graders at 400,000 iterations.

	Accuracy	Avg. Precision	Avg. Recall
one_fc4096	84.51	83.66	80.62
one_fc128	80.14	81.25	72.79
two_fc128	82.11	82.14	76.36

Table 5: Precision, recall, and F1 value of both classes from one_fc4096 model at 400,000 iterations.

	Testing			Training
	Negative	Positive	Average	Average
Precision	85.60	81.73	83.66	99.25
Recall	92.29	68.95	80.62	99.31
F1	88.82	74.80	81.81	99.28

4.4. Scoring Outfits

From the outfit representation Φ , we learn a binary classifier and predict a score. We utilize a multi-layer perceptron (MLP) to learn the mapping function. In this paper, we compare 3 MLPs with various configurations to see the effect of number and size of fully-connected (FC) layers on this problem. The models we used are:

1. one_fc4096: one 4096-d FC layer
2. one_fc128: one 128-d FC layer
3. two_fc128: two 128-d FC layers

Each of fully-connected layers are followed by batch normalization and rectified linear activation (ReLU) with dropout. One 2-d linear layer followed by soft-max activation is added to every models to predict a score. We use multinomial logistic loss to learn the parameters of the grading model.

5. Performance Evaluation

5.1. Evaluation Setup

We learn the grading model from the training split of Polyvore409k dataset, and evaluate the binary classification measures on the testing split. The performance is measured against the ground truth. In this paper, we report the performance of our models without fine-tuning the parameters of the convolutional network for the item feature extraction. We implement the neural network using Caffe framework [23]. We choose cross entropy as a loss function. We train the models for 400,000 iterations using stochastic gradient descent with momentum, where the initial learning rate and momentum are set to 10^{-4} and 0.9, respectively. We measure accuracy, precision, and recall to evaluate the performance. The prediction is counted as correct if it matches the ground truth.



Figure 4: Eight best (top row) and worst (bottom row) outfits judged by our outfit grader

Table 6: Performances of one_fc4096 outfit grader trained by different features: (1) item type, (2) 4-color palette, (3) (1)+(2), (4) ResNet-50 features from grayscale images, (5) ResNet-50 features from RGB images

Feature	Accuracy	Avg. Precision	Avg. Recall
(1) Item types	74.33	71.77	67.02
(2) 4-color palettes	74.53	72.71	66.35
(3) (1)+(2)	78.93	77.57	73.03
(4) ResNet-50 grayscale	81.31	79.35	77.69
(5) ResNet-50 RGB	84.26	83.06	80.73

5.2. Quantitative Results

The accuracy, average precision, and average recall of all models are displayed in table 4. According to the table, one_fc4096, which has 84.51% accuracy, 83.66% average precision, and 80.62% average recall, is clearly the best among the three models. The precision, recall and f1 value of both classes from one_fc4096 model are shown in table 5. Top 8 positive and negative samples from the model are shown in figure 4. Qualitatively, preferred outfits contain items with consistent colors and styles, whereas low-scoring outfits tend to have less common visual elements between items.

From table 5, 92.29% recall for negative class shows that the model is very reliable for pointing out the bad outfit. However, 68.95% for the positive one shows that it tends to judge positive outfit as a negative one as well. When considering that the training performance is almost 100% correct as shown in table 5, we can conclude that the model overfits the training data.

5.3. Color and Item Type Analysis

We conduct another set of experiments to analyze the effect of various features on grading performance. We train one_fc4096 for 100,000 iterations using 5 features: (1) item type, (2) 4-color palette, (3) (1)+(2), (4) ResNet-50 features from grayscale images, and (5) ResNet-50 features from RGB images. Item types are extracted from item name, and 4-color palettes are extracted from item image.

The result in table 6 shows that the item type and color represent the items equally, and the combination of them gives a better representation. However, the composite feature from ResNet-50 outperforms both primitive features, even without the color information. Finally, the color information in the ResNet-50 features affects the performance of outfit grader by 3% classification accuracy.

6. Human Evaluation

Outfit quality is a subjective topic. An outfit that looks chic to one person may look ugly to another. Although an evaluation on testing samples is important, we argue that it might be insufficient to verify the reliability of the approach. We conduct a large-scale human perception evaluation to further assess our model. We use the predictions from one_fc4096 to do evaluations on human perception using Amazon Mechanical Turk (AMT).

6.1. Geographic Trends

People from different regions have different tastes in fashion. Since the one_fc4096 model learned from data from polyvore.com, in order to show that the model successfully learned the compatibilities between fashion items, the model’s predictions should be judged by people from the same region as the training data. After we inspect metadata of all 93,247 outfits that are used as positive samples, we found that they come from 39,590 different users. Around half of them (21,413 users, 54%) did not provide the country. For the remaining 18,177 users, which come from 175 countries, most of them come from United States (8,167 users, 45%), followed by Canada (872 users, 5%), and other countries. As a result, our model’s predictions will be judged by Americans.

6.2. Evaluation Protocol

We setup the experiment as choosing the better outfit from each pair to minimize the effect of absolute bias or personal preference from human subjects. In addition, outfits in each pair must have exactly same outfit parts at the same location in the outfit image, so that only the compatibility of items affects the judgments, not the outfits’ configuration nor number of items in the outfits.

Our hypothesis is, if outfits in the pair have similar quality, people will choose both outfits equally. On the other hand, if the outfit has a large gap in quality, people will definitely choose one over the other. The quality score of each outfit come from our outfit grader. If our outfit grader can reliably judge the quality of the outfit, our hypothesis will be true.

To verify the hypothesis, we select a number of best outfits as references, denoted as *Alpha* (**A**). Then, we select other outfits of different qualities, denoted as *Delta* (**Δ**), and pair them up with best outfits. After that, we show these

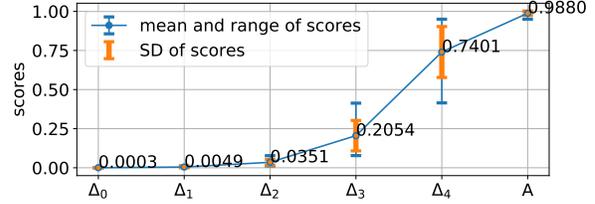


Figure 5: Mean, range, and SD of scores in each samples group Δ and **A**

pairs to human annotators. For each pair, we tell the annotators to choose the better of the two.

Our expectation is that, the difference in quality between outfits in the **A** group and each of Δ group is directly related to the probability that the annotators will choose outfits in the **A** group given the outfits in Δ . Since this experiment is set as pairwise comparisons, we believe that the mentioned probability should be approximated as

$$p(s_\alpha | s_\delta) = s_\alpha / (s_\alpha + s_\delta) \quad (1)$$

where s_α and s_δ are positive probability calculated by our model of outfits in **A** and Δ , respectively.

6.3. Implementation detail

Our outfit’s *score* is the positive probability from the outfit grader. We randomly select 1,000 outfits with the score more than 95% as “**A**” group. After that, we sort outfits with score less than 95% in an ascending order, then divide them into 5 groups, from “ Δ_0 ” which is the group of outfits with the lowest scores, to “ Δ_4 ” which is outfit with the highest scores but still less than 95%.

The experiment consists of 5,000 pairs of outfits. We use outfits from **A** group as “good” outfits, and Δ_j for $j \in \{0, 1, 2, 3, 4\}$ as “bad” outfit. For each $\alpha_i \in \mathbf{A}$, we randomly select an outfit $\delta_{j,i} \in \Delta_j$ for each $j \in \{0, 1, 2, 3, 4\}$ that has exactly same outfit parts as α_i . Our hypothesis is, the visible difference in outfit quality in $(\alpha_i, \delta_{0,i})$ pairs is more than in $(\alpha_i, \delta_{4,i})$. We denote pair $(\alpha_i, \delta_{j,i})$ as $p_{j,i}$ for $j \in \{0, 1, 2, 3, 4\}$ and $i \in \{0, 1, \dots, 999\}$. We show in figure 5 the mean, range, standard deviation of scores in each Δ_j for $j \in \{0, 1, 2, 3, 4\}$ and **A**, and the difference of mean of scores of each group to **A**.

We ask 5 annotators to vote each pair $p_{j,i}$. Each annotator selects the better outfit in each pair, or select “Unable to decide” if the annotator thinks that the outfits looks equally good (or bad). The total number of questions in our experiment equals to: 5 annotators \times 1,000 questions \times 5 $\delta_s = 25,000$ questions. We show examples of outfits in figure 6. For each row, 5 pairs of outfits are created by pairing an outfit in Δ_j for $j \in \{0, 1, 2, 3, 4\}$ with **A**. An example questionnaire is shown in figure 7.



Figure 6: Comparison of outfits used in human evaluation. Each row shows outfits in different quality groups but have the same outfit configuration.

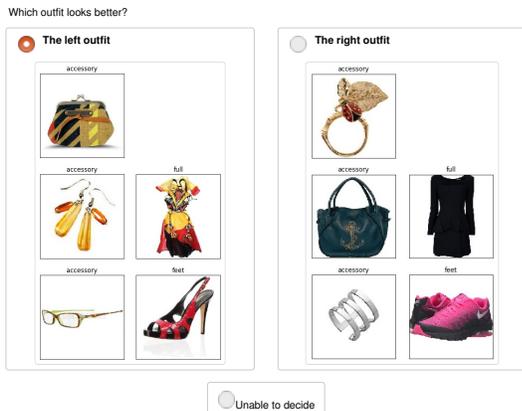


Figure 7: An example of questionnaires used in human evaluation, with associated outfit part of each item

Table 7: Number of “Unable to decide” answers and ties from the experiments comparing outfits in A to Δ_j for $j \in \{0, 1, 2, 3, 4\}$

j	0	1	2	3	4
Number of “Unable to decide” (out of 5,000 questions)	688	820	924	696	422
Number of ties (Out of 1,000 pairs)	63	80	114	111	80

6.4. Evaluation Metrics

We use the term *Matching Ratio* to describe the ratio that human annotators select α_i in pair $(\alpha_i, \delta_{j,i})$ as the better-looking outfit. We also remove the “Unable to decide”

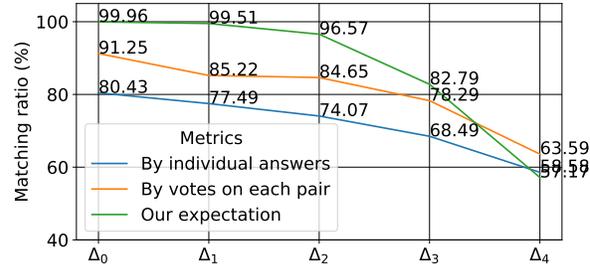


Figure 8: Matching ratio of the prediction to human judgment in each samples group Δ with our expectation

votes, shown in table 7, from the calculation. There are two metrics, matching ratio by individual answer, and by majority vote on each pair. For the latter, we also remove ties, shown in table 7, from the calculation.

6.5. Results

The results, with our expectation as explained in section 6.2, are shown in figure 8. The 91.25% matching ratio by voting shows that the human annotators agree with predictions from our model. Although not perfectly matched, the result has similar trend with our expectation. The result indicates that that the value of our positive probability (score) properly resembles the quality of the outfit.

Regarding the gap between human votes and our model, we have to remind that, the reliability of the human evaluations is not the absolute. As said earlier, fashion is a subjective topic. We might be able to use some small sets of questions to verify the ability of annotators, although this approach introduces absolute bias to the evaluation since people have different tastes in fashion.

7. Application: Outfit Recommendation

If the number of items is not very large, as is often the case with a personal closet, our outfit grader can directly be used as an outfit recommender by generating multiple outfits and ranking them. To demonstrate this usage, we conduct experiments as follows.

7.1. Outfit generation

For generating outfits, we consider four outfit configurations: (1) outer layer with upper- and lower-body, (2) only upper- and lower-body, (3) outer layer with full-body, and (4) full-body only. All configurations include a footwear and at most three optional accessories.

Although it may be reasonable to assume that there are a modest number of clothes, there could be a large number of accessories and generating all possible item combinations is computationally expensive. We test four methods for generating outfits that take efficiency into consideration. The first method (*Ordered Beam Search*) is to regard outfit

generation as a sequence generation problem, and employ a beam search (BS) algorithm. To be specific, in each item step t , the items that belongs to $part_t$ are the possible item extensions, and our “one_fc4096” outfit grader is used as the scoring function. The BS starts from each item in the pool and considers all outfit configurations applicable to the item. It stops when all parts of the outfit are added according to its configuration. We then remove the duplicated outfits and recommend the best outfits based on the score from our outfit grader.

The second method (*Orderless Beam Search*) uses the entire item pool as the possible item extensions at all time steps, while the rest is the same as the first one. The third method (*Partial Beam Search*) generates all possible combination of main parts (outer, upper, lower, full, feet) of the four outfit configurations, from which ten best outfits (based on score from our outfit grader) per outfit configuration per item are kept as base outfits. We then use the beam search to add accessories to those base outfits. The fourth method (*Baseline*) creates 100 outfits per outfit configuration in a random manner. Then, the duplicates are removed and the best ones are recommended.

Each of the four methods outputs 10 best outfits based on the scores from our outfit grader. In the experiments, for all the methods, we set the beam width for beam search to three and include a null item in the item pool as an accessory to give a choice to the beam search to add nothing to an outfit in each “accessory” time steps.

7.2. Evaluation

A good outfit recommender should be able to find sets of well-coordinate items in a pool of apparently random items. To test each recommendation method in terms of this property, we created 957 test cases, each of which contains items from one positive, denoted as P , and two negative samples. These samples are randomly drawn from the testing partition of Polyvore409k dataset. From those items, the recommended outfit, denoted as R , should be similar to the positive samples P . To measure the performance of the recommender, we use four conditions as (1) $P = R$, (2) $P \subset R$ (3) $R \subset P$, (4) $(P = R) \cup (P \subset R) \cup (R \subset P)$. For each method, we regard a recommendation (i.e., top ten recommended outfits) as successful if the condition is met by one of the ten recommended outfits.

Table 8 shows the results. It is seen that *Partial BS* outperforms the baseline in every metrics. The reason why Ordered and Orderless BS perform worse than Partial BS is because our outfit grader is trained using complete outfits, while the early steps of BS rely on the score of partial outfits, which our outfit grader is not trained for. Figure 9 shows successful and unsuccessful recommendations. We argue that the recommended outfits in the failure case are even better than the target positive sample. This is due to



Figure 9: Recommended outfits from Partial BS. Each row shows one test case, where 5 outfits on the right are generated from items in 3 outfits on the left. Outfits with blue border are positive, red are negative, green are *exact match*, cyan are $P \subset R$, and orange are $R \subset P$. The others are recommended outfits that do not meet the conditions.

Table 8: Performance of outfit recommendation by the proposed outfit grader combined with four outfit creation methods. The metrics are (1) $P = R$, (2) $P \subset R$ (3) $R \subset P$, (4) $(P = R) \cup (P \subset R) \cup (R \subset P)$, where P and R denote the positive sample and recommended outfits, respectively.

Approaches	(1)	(2)	(3)	(4)
Ordered BS	11.39	14.11	19.64	32.29
Orderless BS	14.84	20.79	9.40	29.89
Partial BS	34.38	41.80	22.68	59.77
Baseline	8.88	21.53	14.11	36.36

the nature of weakly-supervised data.

8. Conclusion

In this paper, we study outfits as combinations of items by developing outfit graders and outfit recommenders. Given a combination of items as an outfit, our best model can judge if the outfit looks good or not at over 84% accuracy on testing samples, and at 91% matching ratio on evaluations by human annotators. In addition, user can just give a pool of items that user have to our outfit recommender, and it will recommend outfits from the item pool. We also collect a large clothing dataset consisting of over 600,000 clothing items and over 400,000 outfits, and use the dataset to learn and evaluate the outfit graders and recommenders.

Acknowledgements

This work was partly supported by JSPS KAKENHI Grant Number JP15H05919 and CREST, JST Grant Number JPMJCR14D1.

References

- [1] K. Yamaguchi, M. H. Kiapour, L. E. Ortiz, and T. L. Berg, "Parsing clothing in fashion photographs," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3570–3577, IEEE, 2012. 1, 2
- [2] K. Yamaguchi, M. Hadi Kiapour, and T. L. Berg, "Paper doll parsing: Retrieving similar styles to parse clothing items," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 3519–3526, IEEE, 2013. 1, 2
- [3] M. H. Kiapour, K. Yamaguchi, A. C. Berg, and T. L. Berg, "Hipster wars: Discovering elements of fashion styles," in *European Conference on Computer Vision*, pp. 472–488, Springer, 2014. 1
- [4] S. Liu, J. Feng, C. Domokos, H. Xu, J. Huang, Z. Hu, and S. Yan, "Fashion parsing with weak color-category labels," *IEEE Transactions on Multimedia*, vol. 16, no. 1, pp. 253–265, 2014. 1, 2
- [5] E. Simo-Serra, S. Fidler, F. Moreno-Noguer, and R. Urtaun, "A high performance crf model for clothes parsing," in *Asian Conference on Computer Vision*, pp. 64–81, Springer, 2014. 1, 2
- [6] W. Yang, P. Luo, and L. Lin, "Clothing co-parsing by joint image segmentation and labeling," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3182–3189, IEEE, 2014. 1, 2
- [7] S. Vittayakorn, K. Yamaguchi, A. C. Berg, and T. L. Berg, "Runway to realway: Visual analysis of fashion," in *Applications of Computer Vision (WACV), 2015 IEEE Winter Conference on*, pp. 951–958, IEEE, 2015. 1
- [8] A. Veit, B. Kovacs, S. Bell, J. McAuley, K. Bala, and S. Belongie, "Learning visual clothing style with heterogeneous dyadic co-occurrences," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 4642–4650, 2015. 1
- [9] J. Oramas and T. Tuytelaars, "Modeling visual compatibility through hierarchical mid-level elements," *arXiv preprint arXiv:1604.00036*, 2016. 1
- [10] W.-L. Hsiao and K. Grauman, "Learning the latent look: Unsupervised discovery of a style-coherent embedding from fashion images," *arXiv preprint arXiv:1707.03376*, 2017. 1
- [11] S. Liu, J. Feng, Z. Song, T. Zhang, H. Lu, C. Xu, and S. Yan, "Hi, magic closet, tell me what to wear!," in *Proceedings of the 20th ACM International Conference on Multimedia*, pp. 619–628, ACM, 2012. 1, 2, 3
- [12] Y. Hu, X. Yi, and L. S. Davis, "Collaborative fashion recommendation: a functional tensor factorization approach," in *Proceedings of the 23rd ACM International Conference on Multimedia*, pp. 129–138, ACM, 2015. 1, 2, 3
- [13] J. Huang, R. S. Feris, Q. Chen, and S. Yan, "Cross-domain image retrieval with a dual attribute-aware ranking network," in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1062–1070, IEEE, 2015. 1, 2
- [14] Y. Li, L. Cao, J. Zhu, and J. Luo, "Mining fashion outfit composition using an end-to-end deep learning approach on set data," *IEEE Transactions on Multimedia*, 2017. 1, 2, 3
- [15] L.-F. Yu, S. K. Yeung, D. Terzopoulos, and T. F. Chan, "Dressup!: outfit synthesis through automatic optimization," *ACM Trans. Graph.*, vol. 31, no. 6, p. 134, 2012. 1, 2
- [16] E. Simo-Serra, S. Fidler, F. Moreno-Noguer, and R. Urtaun, "Neuroaesthetics in fashion: Modeling the perception of fashionability," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 869–877, IEEE, 2015. 1, 2, 3
- [17] E. Simo-Serra and H. Ishikawa, "Fashion style in 128 floats: joint ranking and classification using weak data for feature extraction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 298–307, IEEE, 2016. 1, 2, 3
- [18] P. Tangseng, Z. Wu, and K. Yamaguchi, "Looking at outfit to parse clothing," *arXiv preprint arXiv:1703.01386*, 2017. 2
- [19] J. McAuley, C. Targett, Q. Shi, and A. Van Den Hengel, "Image-based recommendations on styles and substitutes," in *Proceedings of the 38th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 43–52, ACM, 2015. 2
- [20] Z. Liu, P. Luo, S. Qiu, X. Wang, and X. Tang, "Deepfashion: Powering robust clothes recognition and retrieval with rich annotations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1096–1104, IEEE, 2016. 2
- [21] K. Yamaguchi, T. L. Berg, and L. E. Ortiz, "Chic or social: Visual popularity analysis in online fashion networks," in *Proceedings of the 22nd ACM International Conference on Multimedia*, pp. 773–776, ACM, 2014. 3
- [22] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 770–778, IEEE, 2016. 4
- [23] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell, "Caffe: Convolutional architecture for fast feature embedding," in *Proceedings of the 22nd ACM International Conference on Multimedia*, pp. 675–678, ACM, 2014. 5