

Toward Explainable Fashion Recommendation

Pongsate Tangseng

Takayuki Okatani

Tohoku University

{tangseng, okatani}@vision.is.tohoku.ac.jp

Abstract

Many studies have been conducted so far to build systems for recommending fashion items and outfits. Although they achieve good performances in their respective tasks, most of them cannot explain their judgments to the users, which compromises their usefulness. Toward explainable fashion recommendation, this study proposes a system that is able not only to provide a goodness score for an outfit but also to explain the score by providing reason behind it. For this purpose, we propose a method for quantifying how influential each feature of each item is to the score. Using this influence value, we can identify which item and what feature make the outfit good or bad. We represent the image of each item with a combination of human-interpretable features, and thereby the identification of the most influential item-feature pair gives useful explanation of the output score. To evaluate the performance of this approach, we design an experiment that can be performed without human annotation; we replace a single item-feature pair in an outfit so that the score will decrease, and then we test if the proposed method can detect the replaced item correctly using the above influence values. The experimental results show that the proposed method can accurately detect bad items in outfits lowering their scores.

1. Introduction

Recently, there have been many studies of applying computer vision techniques to various problems of fashion, such as quantifying/measuring goodness of outfits [1–5] and recommending to users outfits from a pool of items [4, 6] or outfits that fit users’s personal preferences [7] or location [8]. However, many of the existing studies, particularly the recent ones that employ CNNs, rely on black-box models, which may provide good performance on respective tasks but cannot explain the reason of their judgments [3, 4, 6, 7]. There are a few attempts to develop models that can provide useful explanations [5, 8], but they require a large amount of manually annotated data for supervised training of the models, which is expensive and usually not



Figure 1: Our system first predicts a goodness score of an input outfit consisting of multiple items. It then identifies which item and what feature is the cause of, for instance, a low score. It is able not only to perform item-level identification (first row) but also to perform feature-level identification (second and third rows).

publicly available.

In this study, we propose a system that is able not only to judge and quantify goodness/badness of an outfit but also to provide a reason(s) of the prediction. Similar to existing methods, our system receives images of multiple items comprising an outfit as inputs and then computes a score quantifying its goodness/badness of the outfit; example inputs are shown in the rows of Fig. 1. This forward computation is done by a part of our system called the outfit grader. To explain the output score, we quantify and use *how large the influence of each item, or of each feature of each item, is on the predicted score*. This enables to identify which item and what feature make the outfit good or bad; examples of the identification are shown in Fig. 1. For this purpose, we represent each item, rigorously its image, with a combination of human-interpretable features, and thereby the identification of the most influential item-feature pair will be a useful explanation of the score.

To measure the influence of item-feature pairs, we em-

ploy the multiplication of an individual feature with the gradient of the output score with respect to the feature. This is similar to the methods for visualizing inference of CNNs, such as the multiplication of an input image with its sensitivity map [9, 10] (i.e., the score gradient with respect to image pixels) and Grad-CAM [11]. The values thus computed are averaged and normalized within each feature of each item to yield our measure of the influence of the item-feature pair, which we call its *Item-Feature Influence Value* (IFIV). Note that our method does not need extra training data other than those for training the outfit grader.

It is usually hard to evaluate explanations provided by AI systems, since their quality can theoretically be evaluated only by humans. Human evaluation is generally costly; moreover, in our case, it is difficult to perform and conveys open problems, as the judgments to be explained are often subjective. To cope with this difficulty, we employ an automatic evaluation method by designing a test for the evaluation that is based on synthesis of datasets. The basic idea is that i) we first replace a single item or its single feature of an outfit so that the resulting score will decrease and ii) we then test if the proposed method can detect the replaced item by identifying the item-feature pair with the maximum IFIV.

The organization of this paper is as follows. We first discuss the related work in Sec. 2. Next, we describe the proposed method for explaining judgments made by our outfit grader on the quality of input outfits in Sec. 3. Section 4 explains and evaluates the outfit grader that is the target of explanation. Experimental results on the proposed method for explaining its judgments are provided in Sec. 5. Section 6 concludes this study.

2. Related Work

2.1. Measuring Goodness of Outfits

There is a growing interest in the application of computer vision techniques to measure the goodness of outfits. The authors of [1] predicted fashionability scores from an outfit image and tags. The authors of [2] use bidirectional LSTM (Bi-LSTM) [12] to learn the compatibility relationship among fashion items by modeling an outfit as a sequence, whereas fully-connected layers are employed in [3, 4]. In [2–4], CNNs trained for generic image recognition are used to extract features for their respective purposes. Overall, the proposed methods in these studies work fairly well for measuring the goodness of outfits, i.e., predicting a score for each outfit. However, these methods lack the ability of providing reasons of the predicted scores.

2.2. Explaining Inference of Models

Recent advances in deep learning have dramatically improved accuracy of many computer vision tasks, such as im-

age classification [13–15], object detection [16], object segmentation [17,18], Visual-Question Answering (VQA) [19–22], etc. These progresses have left behind explanation and understanding of what the deep neural networks have learned as well as how they make inference/judgments. Thus, there is a growing concern particularly about life-critical applications [23]. A number of studies have been conducted to resolve this so far; [11, 24–26] to name a few. LIME [24] is a method for explaining the prediction of a machine learning model for an input, which estimates a linear model that locally approximates the model at the neighborhood of the input, and then uses it for explanation. There are many studies of visualization of inference made by CNNs. The authors of [25] proposed the Class Activation Map (CAM) for a particular class of CNN models, which shows the region in the input image that is responsible for the prediction. This is later extended to Grad-CAM [11], which is applicable to more general CNN models, including image captioning [27–29], and Visual Question Answering (VQA) [19–22].

2.3. Explainable Models for Fashion

The aforementioned computer vision systems for fashion [3, 6, 7, 30, 31] employ black-box models, too, which show fairly good performance for the respective tasks but lack ability of providing reason of inference/judgment. It is not straightforward to apply the above generic methods for explaining machine learning and deep learning models to these systems for fashion, because the problems are basically more complicated (e.g., multiple items contained in an outfit, stratified factors affecting the goodness/badness of an outfit etc.)

There are a few studies that attempt to provide useful explanation on model’s evaluation of outfits [5, 32]. The method proposed in [5] relies on a massive amount of annotated data to train a multi-category attribute predictor and create a composition graph based on pairwise co-occurrence of those predicted attributes in outfits. On the other hand, the method proposed in [32] provides an upper-lower matching recommendation with textual explanation by utilizing comments provided by users of polyvore.com. Although this method does not require manual annotation, it can deal with only two items in each outfit.

3. Explaining Goodness of Outfit

Figure 2 shows an overview of the proposed system. It employs the outfit grader developed in [4], which classifies an input outfit either as positive (a good outfit) or negative (a bad outfit). We wish to explain judgment made by the grader for an outfit, i.e., why it classifies an input outfit as positive or as negative. For this purpose, we evaluate influence of each item and its features on the grader’s judgment. The former (i.e., the influence of each item) provides

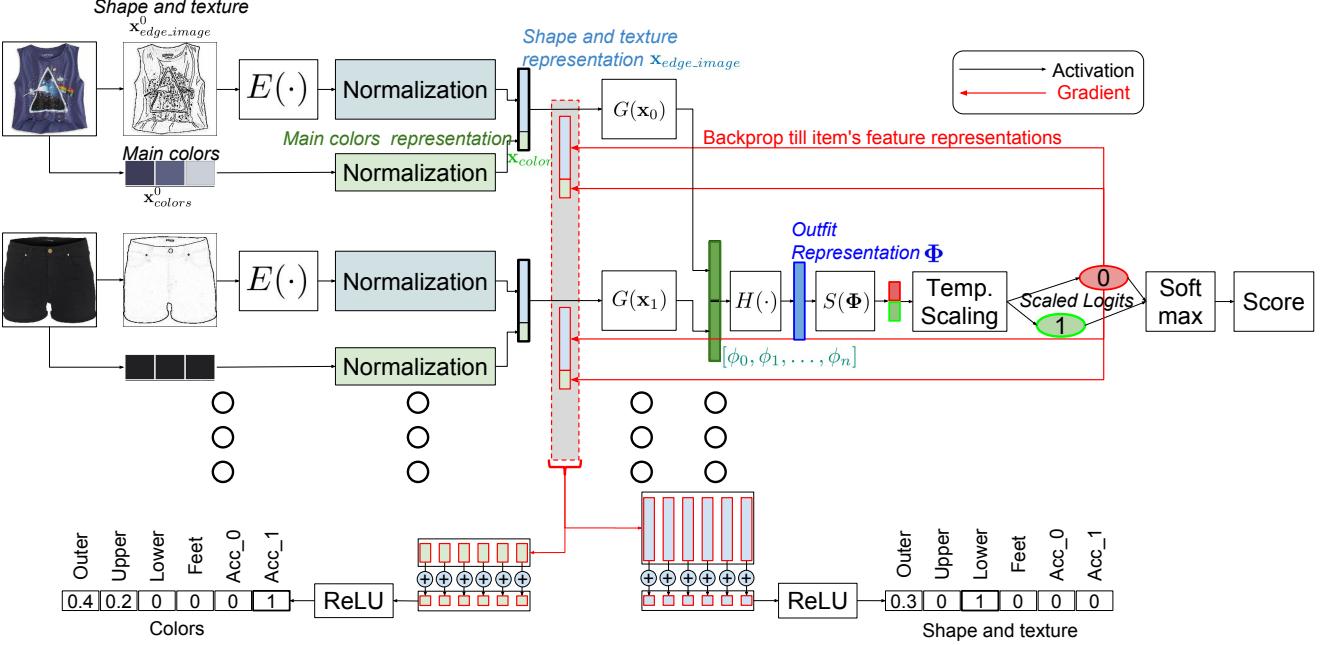


Figure 2: The overview of the proposed system. Given an outfit as a set of items, it extracts *edge_image* and main colors of each item, forward-propagates them through a pretrained CNN, normalization, concatenation, and fully connected layers with ReLU to obtain the score. The system also computes the gradient of the score (rigorously, the logit before softmax) with respect to the representation of each item through backpropagation. The values of gradient at item representation are multiplied with the feature values, and then scaled to $[0, 1]$ range. There is a single value for each item-feature pair. We call the value *Item Feature Influence Value (IFIV)*.

item-level explanations, e.g., *this outfit is bad because of the inclusion of this particular item*. For this, we use the internal features (i.e., penultimate layer activation) that the grader uses. To further enable to obtain deeper explanations, we use human-interpretable features for the purpose, e.g., shape, texture, and colors extracted from the item images comprising the input outfit. To do this, we redesign the grader so that it can make judgments solely from these features.

3.1. Interpretable Item Features

The idea is to represent each item in terms of its *attributes* that are human-interpretable. We also rebuild the grader so that it can judge an input outfit from its attribute representation, and then attempt to explain its judgments according to influence of each attribute on the final score.

There are many candidate for this purpose, such as item type, brand, color, shape, texture, style etc. However, it may be a difficult task even for fashion experts to define such attributes determining the goodness of outfit. Moreover, we also need to be able to accurately predict those attributes from input item images, which will require costly annotation for training a proper model (e.g., a CNN). Additionally, the attributes need to be sufficiently rich so that the grader



Figure 3: Item images with their *edge_image* and main three colors used as their features.

can properly judge goodness of outfits only from them.

Considering these requirements, we choose primitive image features that can be easily extracted from the item images: shape, texture, and colors. To be specific, we first divide contents of item images into color and non-color information. For the former, we extract three dominant colors from each image by finding clusters of pixels in color space. For non-color information, we first convert the image into gray-scale and then extract edges, which are ex-

pected to maintain shape and texture of the item. Figure 3 shows examples of original images, their *edge.image*, and three dominant colors. Their details are given below.

For colors, after removing background from the item image, we apply K-mean clustering [33] to cluster all the pixels in the item image into three main colors in RGB color space. We use their centroids as three dominant colors of the item, yielding a 9-dimensional vector (3 colors \times 3 RGB color values) for each item image, which we will denote by \mathbf{x}_{colors}^0 .

For shape and texture, we extract features in the following way. Let I be the input item image. We first apply the Canny edge detector [34] to I to obtain an edge map I_{e_1} . In parallel, we also apply a simple 3×3 filter f to I as $I_{e_2} = I * f$; f is defined as

$$f = \begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}. \quad (1)$$

We add these two edge-like maps to obtain

$$I_e = I_{e_1} + I_{e_2}. \quad (2)$$

We call its black-white inverted version (i.e., $I_e \leftarrow 255 - clip(I_e, 0, 255)$) *edge.image* of I . We then use a pre-trained convolutional neural network (CNN) to extract an n -dimensional embedding of *edge.image*, which we denote by $\mathbf{x}_{edge.image}^0$, as

$$\mathbf{x}_{edge.image}^0 = E(edge.image), \quad (3)$$

where E is the CNN (rigorously up to its penultimate layer). We will use $\mathbf{x}_{edge.image}^0$ as a representation of shape and texture of an item.

The features \mathbf{x}_{colors}^0 and $\mathbf{x}_{edge.image}^0$ obtained as above are normalized in an element-wise manner over all training samples as

$$x_{i,f} = \frac{x_{i,f}^0 - \mu_{i,f}^0}{\sigma_{i,f}^0}, \quad f \in \{edge.image, colors\}, \quad (4)$$

where i is the index of the vectors; $\mu_{i,f}^0$ and $\sigma_{i,f}^0$ are the mean and standard deviation of the i -th element of \mathbf{x}_f^0 . The identical standardization is applied to the features of an input outfit at test time.

Finally, we concatenate n -dimensional $\mathbf{x}_{edge.image}$ and 9-dimensional \mathbf{x}_{colors} and denote the resultant $n + 9$ -dimensional vector by \mathbf{x} .

$$\mathbf{x} = [\mathbf{x}_{edge.image}^\top, \mathbf{x}_{colors}^\top]^\top, \quad (5)$$

which gives a representation of an item. An outfit consists of multiple items, each of which occupies a specific outfit part. We will denote the representation of an item of the i -th part by \mathbf{x}_i .

3.2. Outfit Grader

Our outfit grader is basically the same as the one proposed in [4] except the representation of items described above. We summarize its design here. The input is an outfit consisting of n items, each of which occupies a different part. Given the feature of each item as mentioned above, our grader first transforms it by a trainable item encoder G as

$$\phi_i = G(\mathbf{x}_i). \quad (6)$$

In our experiments, we use a few fully-connected layer for G . The representations of n items are then concatenated and transformed to the representation Φ of the entire outfit as

$$\Phi = H([\phi_0, \phi_1, \dots, \phi_n]), \quad (7)$$

where H is a trainable outfit encoder, for which we employ a single fully-connected layer (followed by BN and ReLU).

The grader performs binary classification on the representation Φ of the input outfit O . To do this, the outfit representation is transformed by a single fully-connected layer S to two logits $\mathbf{s} = [s_{pos}, s_{neg}]$ as $\mathbf{s} = S(\Phi)$. Then they are normalized by softmax to yield scores for positive and negative classifications. Denoting the score for O being positive by $F(O)$, it is given by

$$F(O) = \sigma_{pos}(\mathbf{s}) = \frac{\exp(s_{pos})}{\exp(s_{pos}) + \exp(s_{neg})}. \quad (8)$$

For the CNNs extracting item features (e.g., $\mathbf{x}_{edge.image}$), we use those pretrained on other tasks such as object recognition. Thus, the learnable parameters in the grader are in G , H , and S . They are learned by minimizing a cross-entropy loss on training data consisting of pairs of outfit O and the ground-truth label (i.e., positive or negative).

Calibration of Outfit Scores It is known [35] that modern deep neural networks employing softmax for multi-class classification tend to be over-confident, that is, the score of the predicted class, or *confident* (i.e., the max of softmax outputs), tends to be large and even close to one, even if the prediction is wrong. We found that this is exactly the case with our implementation of the outfit grader [4]. A simple but effective method to alleviate this overconfidence is to perform calibration of the softmax outputs using temperature scaling [35, 36]. To be specific, we replace \mathbf{s} in the softmax (8) with \mathbf{s}/T . T is determined using validation samples so that the resulting score $F(O)$ is as close to classification accuracy as possible; then the score will better represent confidence of the prediction. We use $\hat{q} = 100 \cdot F(O)$ (in percent) as the fashionability score of an outfit O .

3.3. Item Feature Influence Value (IFIV)

Suppose that we input an outfit to the above grader and receive its judgment. To explain the judgment, we evaluate influence of each feature of each item. If the judgment is negative and a particular feature of an item has large influence on it, we regard that feature of the item to be the reason for the negativity; the same is true for a positive judgment.

To be specific, we define the influence on the logit s_c ($c \in \{neg, pos\}$) of a feature f ($\in \{edge_image, colors\}$) of i -th item, denoted by $\mathbf{x}_{i,f}$, as follows. We first compute

$$\mathbf{g}_{i,f} = \mathbf{x}_{i,f} \odot \frac{\partial s_c}{\partial \mathbf{x}_{i,f}}, \quad (9)$$

where \odot is element-wise multiplication. Note that the logit s_c here is the temperature-scaled version mentioned above. A similar method is used for visualization of CNNs for object classification, where the pixel-wise multiplication of an input image and the gradient of a class score with respect to its pixels is used to show which part positively or negatively affects the score and which part has no influence on it. As we consider influence of only each feature, not its element, we compute the sum over its all elements as

$$v_{i,f} = \text{ReLU} \left(\sum_k g_{i,f,k} \right), \quad (10)$$

where $g_{i,f,k}$ is the k -th element of $\mathbf{g}_{i,f}$. To enable comparison between different features, we normalize them as

$$IFIV_{i,f} = \frac{v_{i,f}}{\max_j(v_{j,f})} \quad (11)$$

Figure 2 shows the diagram explaining how *Item Feature Influence Value (IFIV)* of each item feature is computed.

4. Evaluation of the Outfit Grader

4.1. Precision Accuracy vs. Interpretability

We redesign the outfit grader for the purpose of improved explanability. The original model [4] is designed to be an end-to-end model receiving raw item images as inputs, aiming at the best prediction accuracy of outfit quality. Our redesigned model receives hand-engineered features extracted from item images for the sake of explanability. This will sacrifice accuracy of outfit quality prediction. We conducted experiments to examine this.

Model architecture We compare two models that differ only in the item representation \mathbf{x} . One is the model we described in Sec. 3. The other is a baseline model, which uses a CNN feature directly extracted from RGB item images; to be specific, the feature of each item is given by $\mathbf{x} = E(RGB_image)$, where E is a pretrained CNN that is

Table 1: Training, validation, and testing accuracy and average f1 of two outfit graders (a baseline and the interpretable model) on Polyvore409k dataset [4].

Partition	Metric	Model	
		Baseline	Interpretable
Train	Acc.	98.41	99.62
	Avg. F1	98.20	99.57
Validation	Acc.	83.19	80.78
	Avg. F1	81.86	79.48
Test	Acc.	79.19	76.69
	Avg. F1	74.11	71.11

the same as the one used to extract $\mathbf{x}_{edge_image}^0$. The configurations and parameters that are shared by the two models are as follows:

- For the feature extractor E , we employ ImageNet-pretrained InceptionV3 [15]. The activation of $pool5$ layer for an input item image is used for \mathbf{x} , which forms a 2048-dimensional vector.
- An identity function is used for the item encoder G .
- A single fully-connected layer with 4096 units is used for the outfit encoder H , followed by batch normalization [37] and ReLU [38] activation function.
- The both models are trained for 30 epochs with learning rate $1e - 4$ and batch size 100 on Polyvore409k dataset [4].

Results Table 1 shows the results. Accuracy indicates that of binary classification, where a prediction is considered to be correct if it matches the ground truth. As expected, the baseline model shows better performance than the interpretable model by 2.50% accuracy and 3.00% average f1. This is a noticeable gap but is arguably not so large to make the explanation by the interpretable model meaningless.

Configuration of Outfit Grader To recover the performance drop as much as possible and further achieve better prediction accuracy, we tested a number of configurations of the interpretable grader. To be specific, we tested different configurations of the item encoder G and the outfit encoder H . The configurations and their performance on testing samples are shown in Table 2. Since the model #4 has the best performance, we will use this model for the experiments on explainability using feature influence values. Figure 4 shows examples of judgments of the grader; outfits with the highest score and those with the lowest scores.



Figure 4: The best (upper) and worst (lower) eight outfits from testing partition of Polyvore409k dataset according to our outfit grader.

Table 2: Testing accuracy and average f1 of various configurations of outfit grader after training for 30 epochs of Polyvore409k dataset [4]. Each cell in the Item Encoder G and Outfit Encoder H column specify the size of FC layer in the FC block. The \times indicates multiple FC blocks.

#	Item Encoder G	Outfit Encoder H	Acc.	Avg. F1
1	-	4096	76.69	71.11
2	1024	4096	78.93	73.43
3	512×256	4096	79.34	75.19
4	512×256	2048	79.45	75.76
5	256×128	4096	79.00	74.47
6	256×128	2048	78.78	74.70

4.2. Effect of Calibration of Score (Confidence)

As mentioned in Sec. 3.2, we employ the temperature scaling to calibrate the outfit score (or confidence) \hat{q} . Figure 5 shows the reliability diagrams [39, 40] before and after the calibration. Searching for the best value for the temperature T on the validation samples yielded $T = 6.97$. To do this, we split all the testing samples into 10 bins with an equal width, using which we plot the expected accuracy of samples in each bin against the average confidence from the outfit scores. A perfectly calibrated model will yield an identity relation between them. We also calculated expected calibration error (ECE) [41], the difference in expectation between confidence and accuracy. ECE is reduced from 13.90 and 16.02 before the calibration to 1.10 and 2.16 after calibration for validation and testing partition of Polyvore409k dataset [4] respectively. Figure 6 shows distributions of outfit scores for samples with positive labels and those with negative labels. The distributions with the temperature scaling clearly have a much wider spread, making the score more meaningful. We can conclude from Figs. 5 and 6 that the temperature scaling is able to calibrate the outfit scores.

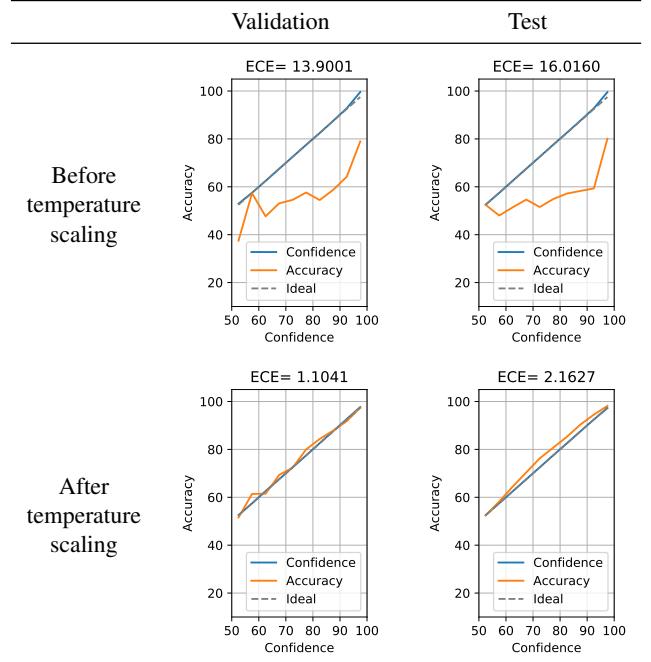


Figure 5: Reliability diagrams and ECE values before and after temperature scaling for validation and testing partition of Polyvore409k dataset [4]. Confidence is equivalent to the outfit score.

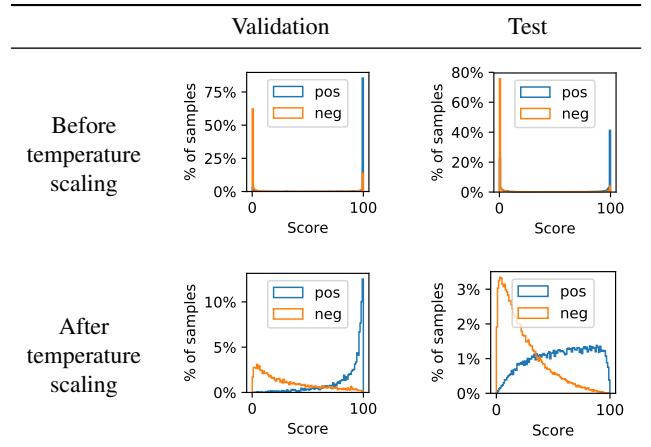


Figure 6: Distribution of outfit scores before and after temperature scaling for positive and negative samples in validation and testing partition of Polyvore409k dataset [4].

5. Experimental Results

We conducted experiments to evaluate the proposed method for explaining judgment of the outfit grader. For the grader, we used the $512 \times 256\text{-}2048$ outfit grader from Table 2.

5.1. Experimental Design

Suppose that an outfit is bad (i.e., not fashionable) due to a single item contained in it. There should also be a reason why the item does not match the outfit and makes it bad, e.g., because of its incompatible color or its unmatched shape and texture. We want to identify the item as well as the reason for the bad outfit.

Based on the proposed framework, this is formulated as a task of identifying the item-feature pair that has the most negative influence on an input outfit. We apply the proposed method to this task and evaluate its performance.

For this purpose, we create a set of negative outfits from positive ones in the dataset in the following way. For a positive outfit, we choose an item from those contained in it and then replace its feature $f \in \{\text{edge_image}, \text{colors}\}$ with that of other items belonging to the same outfit part. We also ensure that the replacement does decrease the outfit score. Note that we are interested here not in the correctness of the judgment of the outfit grader but in how well its judgment can be explained, more precisely, accuracy of the proposed method identifying the item-feature pair lowering the score.

Detailed procedures for the creation of data are as follows:

1. 1,000 base outfits with the highest scores are chosen from the test partition of Polyvore409k dataset [4]. Their average score is 97.16 (out of 100).
2. For each item and its feature f in each base outfit, we create 10 mod samples in the following way:
 - 2.1 500 mod samples are first created by replacing the item-feature f in the base sample with that of a randomly chosen item of the same outfit part from the test partition of the dataset.
 - 2.2 Their scores are computed by the outfit grader and the worst ten samples are selected and all the others are discarded.

Step 2.2 ensures that the grader gives low scores to the created outfits with a replaced item-feature pair. For the two features of *edge_image* and *colors*, the above procedure produces two datasets, which we call *edge_image-wise* and *colors-wise* samples, respectively. Two examples of created negative samples are shown in Fig. 8. Additionally, we create “item-wise” samples by replacing the entire item in Step 2.1. The statistics of the base samples and the three types of negative samples are shown in Table 3. The distributions of scores for these samples are shown in Fig. 7.

5.2. Results

We apply our method to the three types of samples created as explained above. To be specific, inputting each sample to the grader, which yield a lower score as explained

Table 3: Statistics of the base samples and the negative samples created from them. The three types of negative samples, i.e., *edge_image*-wise, *colors*-wise, and item-wise, have identical statistics by their construction.

Sample type	Number of samples containing following		
	outfit parts	number of items	
Base sample	Outer	385	3 items
	Upper	615	4 items
	Lower	664	5 items
	Full	379	6 items
	Feet	953	7 items
	Accessory0	977	8 items
	Accessory1	903	1
	Accessory2	722	Total 1,000
Outfit flaw detection sample	Outer	3,850	3 items
	Upper	6,150	4 items
	Lower	6,640	5 items
	Full	3,790	6 items
	Feet	9,530	7 items
	Accessory0	9,770	8 items
	Accessory1	9,030	1
	Accessory2	7,220	Total 55,980

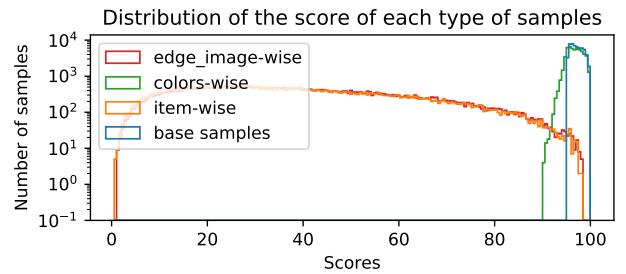


Figure 7: The distribution of scores of each type of samples.

above, we compute *IFIVs* for the score defined in (11). We then find the part with the minimum *IFIV*, or equivalently, that the maximum negative *IFIV*, for each f , as

$$i^* = \arg \max_i (-\text{IFIV}_{i,f}). \quad (12)$$

We regard the prediction i^* as correct if it matches the true item, which is the replaced one when creating the negative sample. Figure 8 shows examples of *IFIVs* for different types of samples. It is seen that the replaced item-feature pairs yield high negative *IFIVs*, meaning that our method can successfully detect the item lowering the outfit score with the reason why it is bad (i.e., the feature lowering the outfit score).

Table 4 show the performance over all the samples. The proposed method can detect the replaced items for *item-wise* samples with 99.52% accuracy and those for

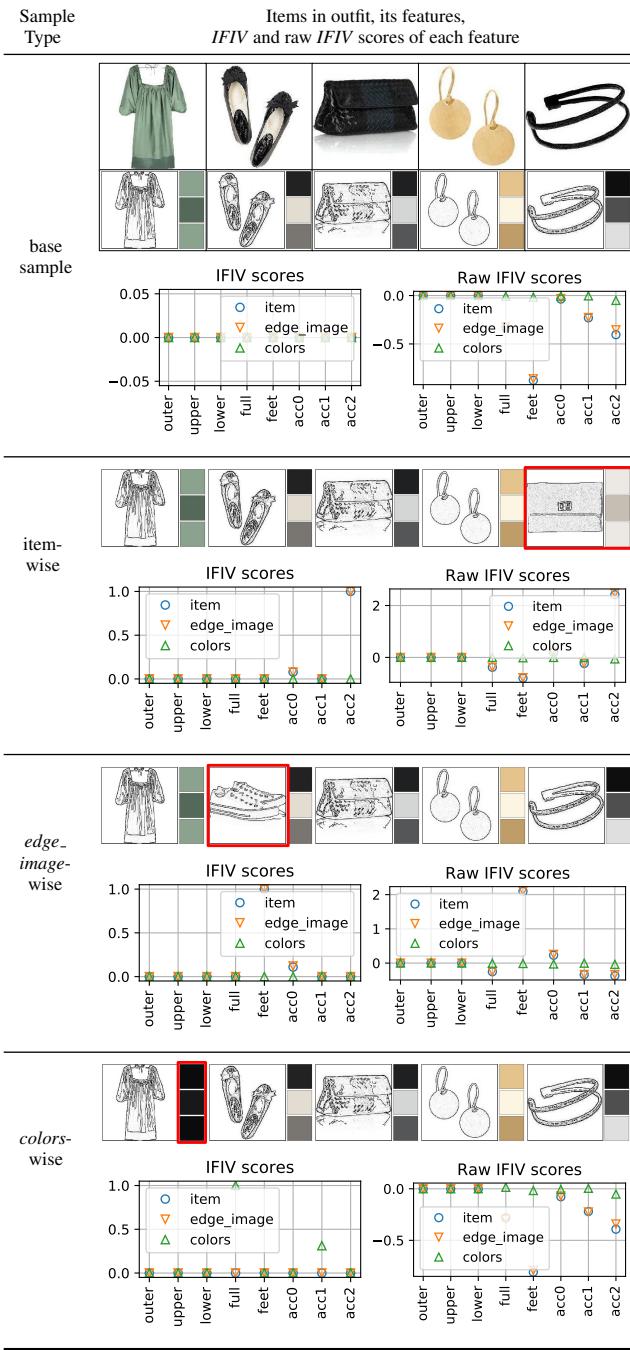


Figure 8: Examples of IFIVs computed by the proposed method. The red boxes indicate the replaced entities from the original high-quality outfits, which makes the new outfits have low outfit scores. IFIV scores mean negative IFIV values.

edge-image-wise samples with 99.49% accuracy, respectively. The accuracy for *colors-wise* samples is 85.37% and lower. This is due to the fact that the scores of the *colors*-

Table 4: Overall accuracy (%) of detection of replaced item-feature pairs.

Method	sample type	prediction accuracy
Random	all	17.86
Proposed method	item-wise	99.52
	<i>edge_image</i> -wise	99.49
	<i>colors</i> -wise	85.37

Table 5: Accuracy (%) of replaced item-feature detection for different numbers of items contained in each outfit.

Number of items	By chance	Proposed method (by sample type)		
		item	<i>edge_image</i>	<i>colors</i>
3	33.33	99.78	99.33	93.78
4	25.00	98.94	98.75	90.09
5	20.00	99.59	99.54	89.34
6	16.67	99.67	99.66	83.52
7	14.29	99.42	99.47	81.31
8	12.50	75.00	75.00	52.50

Table 6: Accuracy (%) of replaced item-feature detection classified by different outfit parts. Note that there are eight outfit parts in Polyvore409k dataset.

Outfit part	By chance	Proposed method (by sample type)		
		item	<i>edge_image</i>	<i>colors</i>
outer	15.44	99.25	99.17	78.91
upper	16.97	99.11	99.32	79.93
lower	16.97	99.94	99.89	86.07
full	18.77	95.57	95.09	85.65
feet	17.36	100.00	100.00	93.05
accessory0	17.31	99.90	99.84	86.82
accessory1	16.96	99.99	99.99	85.71
accessory2	16.50	99.96	100.00	80.11

wise samples tend to be higher and their gap to the original outfits are smaller than the other two types, as shown in Fig. 7. That said, its accuracy is fairly good considering the chance rate of 17.86%.

Table 5 shows accuracy values for different numbers of items. They are quite consistence for *item*- and *edge_image*-wise samples, except for the outfit with eight items. Note that there is only one out of 1,000 base samples that has eight items, as shown in Table 3, and thus the performance for eight items could be statistically unreliable. For *colors*-wise samples, there is a tendency that the accuracy decreases as the number of items increases.

Table 6 shows accuracy values calculated for each part of outfits. It is seen that for *item*- and *edge_image*-wise samples, the performance are almost the same across all outfit

parts, except the full outfit part showing slightly lower accuracy. For *colors*-wise samples, the accuracies are lower than the other two types and are somewhat different for different parts.

6. Conclusion

In this paper, we have proposed a novel method for item-feature-wise explanation of outfits. The method can quantify the effect of interpretable features of each item on the goodness of an outfit with the proposed *Item Feature Influence Value (IFIV)*. It does not need any item-level attribute annotation. Using the *IFIV* of each item-feature pair in an outfit, we can detect the bad item in an outfit lowering its score by finding the item-feature pair with the maximum negative *IFIV*. The experiments have shown that our method can detect the bad items at 99.52, 99.49, and 85.37%, for datasets of item-wise, *edge_image*-wise, and *colors*-wise samples, respectively.

References

- [1] E. Simo-serra, S. Fidler, F. Moreno-noguer, R. Urtau, and I. D. Rob, “Neuroaesthetics in Fashion : Modeling the Perception of Fashionability,” *Proceedings of The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. [1](#) [2](#)
- [2] X. Han, Z. Wu, Y.-G. Jiang, and L. S. Davis, “Learning fashion compatibility with bidirectional lstms,” in *Proceedings of the 2017 ACM on Multimedia Conference*, pp. 1078–1086, ACM, 2017. [1](#) [2](#)
- [3] 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](#)
- [4] P. Tangseng, K. Yamaguchi, and T. Okatani, “Recommending outfits from personal closet,” in *Proceedings of IEEE Winter Conference on Applications of Computer Vision (WACV)*, vol. 00, pp. 269–277, Mar 2018. [1](#) [2](#) [4](#) [5](#) [6](#) [7](#)
- [5] Z. Feng, Z. Yu, Y. Yang, Y. Jing, J. Jiang, and M. Song, “Interpretable partitioned embedding for customized multi-item fashion outfit composition,” in *Proceedings of the 2018 ACM on International Conference on Multimedia Retrieval*, pp. 143–151, ACM, 2018. [1](#) [2](#)
- [6] 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](#)
- [7] 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](#)
- [8] X. Zhang, J. Jia, K. Gao, Y. Zhang, D. Zhang, J. Li, and Q. Tian, “Trip outfits advisor: Location-oriented clothing recommendation,” *IEEE Transactions on Multimedia*, vol. 19, no. 11, pp. 2533–2544, 2017. [1](#)
- [9] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus, “Intriguing properties of neural networks,” *arXiv preprint arXiv:1312.6199*, 2013. [2](#)
- [10] D. Smilkov, N. Thorat, B. Kim, F. Vigas, and M. Wattberg, “Smoothgrad: removing noise by adding noise,” *arXiv preprint arXiv:1706.03825*, 2017. [2](#)
- [11] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, *et al.*, “Grad-cam: Visual explanations from deep networks via gradient-based localization,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pp. 618–626, 2017. [2](#)
- [12] A. Graves, “Supervised sequence labelling,” in *Supervised sequence labelling with recurrent neural networks*, pp. 5–13, Springer, 2012. [2](#)
- [13] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *CoRR*, vol. abs/1409.1556, 2014. [2](#)
- [14] 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. [2](#)
- [15] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna, “Rethinking the inception architecture for computer vision,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2818–2826, 2016. [2](#) [5](#)
- [16] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, pp. 91–99, 2015. [2](#)
- [17] J. Long, E. Shelhamer, and T. Darrell, “Fully convolutional networks for semantic segmentation,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3431–3440, IEEE, 2015. [2](#)

- [18] J. Dai, K. He, and J. Sun, “Instance-aware semantic segmentation via multi-task network cascades,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3150–3158, IEEE, 2016. [2](#)
- [19] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh, “Vqa: Visual question answering,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 2425–2433, 2015. [2](#)
- [20] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu, “Are you talking to a machine? dataset and methods for multilingual image question,” in *Advances in neural information processing systems*, pp. 2296–2304, 2015. [2](#)
- [21] M. Malinowski, M. Rohrbach, and M. Fritz, “Ask your neurons: A neural-based approach to answering questions about images,” in *Proceedings of the IEEE International Conference on Computer Vision*, pp. 1–9, 2015. [2](#)
- [22] M. Ren, R. Kiros, and R. Zemel, “Exploring models and data for image question answering,” in *Advances in neural information processing systems*, pp. 2953–2961, 2015. [2](#)
- [23] Z. C. Lipton, “The mythos of model interpretability,” *arXiv preprint arXiv:1606.03490*, 2016. [2](#)
- [24] M. T. Ribeiro, S. Singh, and C. Guestrin, “Why should i trust you?: Explaining the predictions of any classifier,” in *Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining*, pp. 1135–1144, ACM, 2016. [2](#)
- [25] B. Zhou, A. Khosla, A. Lapedriza, A. Oliva, and A. Torralba, “Learning deep features for discriminative localization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 2921–2929, 2016. [2](#)
- [26] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba, “Network dissection: Quantifying interpretability of deep visual representations,” *arXiv preprint arXiv:1704.05796*, 2017. [2](#)
- [27] X. Chen, H. Fang, T.-Y. Lin, R. Vedantam, S. Gupta, P. Dollár, and C. L. Zitnick, “Microsoft coco captions: Data collection and evaluation server,” *arXiv preprint arXiv:1504.00325*, 2015. [2](#)
- [28] J. Johnson, A. Karpathy, and L. Fei-Fei, “Densecap: Fully convolutional localization networks for dense captioning,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 4565–4574, 2016. [2](#)
- [29] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 3156–3164, 2015. [2](#)
- [30] 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](#)
- [31] K. Matzen, K. Bala, and N. Snavely, “StreetStyle: Exploring world-wide clothing styles from millions of photos,” *arXiv preprint arXiv:1706.01869*, 2017. [2](#)
- [32] Y. Lin, P. Ren, Z. Chen, Z. Ren, J. Ma, and M. de Rijke, “Explainable fashion recommendation with joint outfit matching and comment generation,” *arXiv preprint arXiv:1806.08977*, 2018. [2](#)
- [33] S. Lloyd, “Least squares quantization in pcm,” *IEEE transactions on information theory*, vol. 28, no. 2, pp. 129–137, 1982. [4](#)
- [34] J. Canny, “A computational approach to edge detection,” *IEEE Transactions on pattern analysis and machine intelligence*, no. 6, pp. 679–698, 1986. [4](#)
- [35] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, “On calibration of modern neural networks,” *arXiv preprint arXiv:1706.04599*, 2017. [4](#)
- [36] J. Platt *et al.*, “Probabilistic outputs for support vector machines and comparisons to regularized likelihood methods,” *Advances in large margin classifiers*, vol. 10, no. 3, pp. 61–74, 1999. [4](#)
- [37] S. Ioffe and C. Szegedy, “Batch normalization: Accelerating deep network training by reducing internal covariate shift,” *arXiv preprint arXiv:1502.03167*, 2015. [5](#)
- [38] V. Nair and G. E. Hinton, “Rectified linear units improve restricted boltzmann machines,” in *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814, 2010. [5](#)
- [39] M. H. DeGroot and S. E. Fienberg, “The comparison and evaluation of forecasters,” *The statistician*, pp. 12–22, 1983. [6](#)

- [40] A. Niculescu-Mizil and R. Caruana, “Predicting good probabilities with supervised learning,” in *Proceedings of the 22nd international conference on Machine learning*, pp. 625–632, ACM, 2005. [6](#)
- [41] M. P. Naeini, G. F. Cooper, and M. Hauskrecht, “Obtaining well calibrated probabilities using bayesian binning..,” in *AAAI*, pp. 2901–2907, 2015. [6](#)