



Attention-based Fusion for Outfit Recommendation

Katrien Laenen, Marie-Francine Moens

University: KU Leuven

Motivation

- Explosive growth of e-commerce content on the Web
- Recommendation systems are essential to overcome consumer overchoice
- Limited support for users looking for a full outfit



Frequently Bought Together



Reiss
EUR 284.18



New!
Sanctuary
EUR 56.84

2



Madewell
EUR 86.22

LOOKS



EUR 13.87



EUR 32.75

EUR 161.84

...



SOLD OUT



EUR 238.91



EUR 934.43

...



EUR 260.05

...

Problem definition

Outfit
compatibility
prediction



Outfit
completion



+



Challenges

- **Item understanding**
 - Capture important fine-grained product features in the item representation
 - Effectively fuse the information in the product image and description
→ attention-based fusion
- **Item matching**
 - Compatibility is a complex relationship (e.g., not transitive)



applied embellished fit flare dress



pre-owned christian louboutin confusalta T-strap platform peep toe pum

Overview

- Motivation
- Problem definition
- Challenges
- Methodology
 - Baseline model: Common space fusion
 - Our model: Attention-based fusion
- Datasets
- Experiments
- Results
- Conclusions

Methodology: Common space fusion

Baseline model

Common space fusion method of Vasileva et al. (2018)

Input

A triplet of image embeddings $(\mathbf{x}_{(u)}, \mathbf{x}_{(v)}^+, \mathbf{x}_{(v)}^-)$ and a triplet of corresponding description embeddings $(\mathbf{t}_{(u)}, \mathbf{t}_{(v)}^+, \mathbf{t}_{(v)}^-)$

Example: u = Dresses

v = Shoes

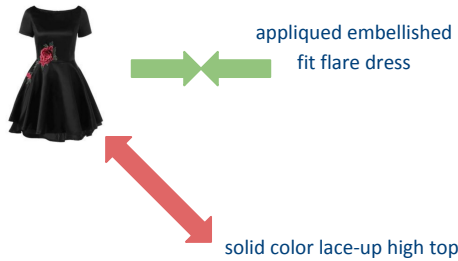
$(\mathbf{x}_{(u)}, \mathbf{x}_{(v)}^+, \mathbf{x}_{(v)}^-) =$ 

$(\mathbf{t}_{(u)}, \mathbf{t}_{(v)}^+, \mathbf{t}_{(v)}^-) =$ (applied embellished fit flare dress 6, pre-owned christian louboutin confusalta T-strap platform peep toe pum, solid color lace-up high top)

Methodology: Common space fusion

Multimodal semantic space

visual-semantic loss



$$\mathcal{L}_{vse} = \frac{\mathcal{L}_{vse, \mathbf{x}(u)} + \mathcal{L}_{vse, \mathbf{x}(v)^+} + \mathcal{L}_{vse, \mathbf{x}(v)^-}}{3}$$

$$\mathcal{L}_{vse, \mathbf{x}(u)} = \frac{\ell(W_i \mathbf{x}(u), W_s \mathbf{t}(u), W_s \mathbf{t}(v)^+) + \ell(W_i \mathbf{x}(u), W_s \mathbf{t}(u), W_s \mathbf{t}(v)^-)}{2}$$

with $\ell(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \max(0, f(\mathbf{x}, \mathbf{z}) - f(\mathbf{x}, \mathbf{y}) + m)$

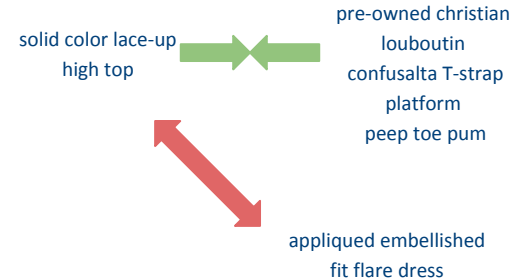
$$\text{and } f(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \cdot \mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|}$$

visual similarity loss



$$\mathcal{L}_{vsim} = \frac{\ell(W_i \mathbf{x}(v)^+, W_i \mathbf{x}(v)^-, W_i \mathbf{x}(u)) + \ell(W_i \mathbf{x}(v)^-, W_i \mathbf{x}(v)^+, W_i \mathbf{x}(u))}{2}$$

textual similarity loss



$$\mathcal{L}_{tsim} = \frac{\ell(W_s \mathbf{t}(v)^+, W_s \mathbf{t}(v)^-, W_s \mathbf{t}(u)) + \ell(W_s \mathbf{t}(v)^-, W_s \mathbf{t}(v)^+, W_s \mathbf{t}(u))}{2}$$

Methodology: Common space fusion

Type-specific compatibility spaces

compatibility loss

$$\mathcal{L}_{comp} = \ell(W_c^{(u,v)} W_i \mathbf{x}_{(u)}, W_c^{(u,v)} W_i \mathbf{x}_{(v)}^+, W_c^{(u,v)} W_i \mathbf{x}_{(v)}^-)$$

with $\ell(\mathbf{x}, \mathbf{y}, \mathbf{z}) = \max(0, f(\mathbf{x}, \mathbf{z}) - f(\mathbf{x}, \mathbf{y}) + m)$

$$\text{and } f(\mathbf{x}, \mathbf{y}) = \frac{\mathbf{x}^T \cdot \mathbf{y}}{\|\mathbf{x}\| \cdot \|\mathbf{y}\|}$$

Compatibility space for *Dresses* and *Shoes*:



Training

complete loss

$$\mathcal{L} = \mathcal{L}_{comp} + \lambda_1 \mathcal{L}_{vsim} + \lambda_2 \mathcal{L}_{tsim} + \lambda_3 \mathcal{L}_{vse}$$

Methodology: Attention-based fusion

Input

A triplet of **region-level** image feature: $(\mathbf{x}_{1:N(u)}, \mathbf{x}_{1:N(v)}^+, \mathbf{x}_{1:N(v)}^-)$ and a triplet of corresponding description-level features $(\mathbf{t}_{(u)}, \mathbf{t}_{(v)}^+, \mathbf{t}_{(v)}^-)$ **or** word-level features $(\mathbf{t}_{1:M(u)}, \mathbf{t}_{1:M(v)}^+, \mathbf{t}_{1:M(v)}^-)$ (depends on attention mechanism)

Multimodal semantic space

Average region-level and word-level representations, i.e., $\mathbf{x}_{(u)} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_{i(u)}$, to compute losses
Use **attention** to fuse the visual and textual information to obtain a triplet of multimodal item representations $(\mathbf{m}_{(u)}, \mathbf{m}_{(v)}^+, \mathbf{m}_{(v)}^-)$

Type-specific compatibility spaces

$$\mathcal{L}_{comp} = \ell(W_c^{(u,v)} \frac{\mathbf{m}_{(u)}}{W_i \mathbf{x}_{(u)}}, W_c^{(u,v)} \frac{\mathbf{m}_{(v)}^+}{W_i \mathbf{x}_{(v)}^+}, W_c^{(u,v)} \frac{\mathbf{m}_{(v)}^-}{W_i \mathbf{x}_{(v)}^-})$$

Methodology: Visual dot product attention

Input

Region-level image features $X \in \mathbb{R}^{N \times d_g}$

Description-level text feature $t \in \mathbb{R}^{d_g}$

Visual attention weights and context vector

$$a_i = \tanh(\mathbf{x}_i) \cdot \tanh(\mathbf{t})$$

$$\mathbf{c} = \sum_{i=1}^N \alpha_i \mathbf{x}_i, \text{ with } \alpha_i = \text{softmax}([a_1, a_2, \dots, a_N])_i$$

Multimodal item representation

$$[\mathbf{c}; \mathbf{t}]$$

Methodology: Stacked visual attention

Input

Region-level image features $X \in \mathbb{R}^{N \times d_g}$

Description-level text feature $t \in \mathbb{R}^{d_g}$

Visual attention weights and context vector

Computed in R attention hops

$$\mathbf{a}^{(r)} = \mathbf{w}_p^{(r)} \tanh(W_v^{(r)} X^T \oplus (W_t^{(r)} \mathbf{q}^{(r-1)} + \mathbf{b}_s^{(r)}))$$

$$\mathbf{c}^{(r)} = \boldsymbol{\alpha}^{(r)} X, \text{ with } \boldsymbol{\alpha}^{(r)} = \text{softmax}(\mathbf{a}^{(r)})$$

$$\mathbf{q}^{(r)} = \mathbf{q}^{(r-1)} + \mathbf{c}^{(r)}$$

Multimodal item representation

$$[\mathbf{q}^{(R)}; t]$$

11

Methodology: Co-attention

Input

Region-level image features $X \in \mathbb{R}^{N \times d_g}$

Word-level text features $Y \in \mathbb{R}^{M \times d_g}$

Textual attention

$$\mathbf{a}^t = \text{Convolution1D}_{t,2}(\text{ReLU}(\text{Convolution1D}_{t,1}(Y)))$$

$in=d_g, out=1, k=1$ $in=d_g, out=d_g, k=1$

$$\mathbf{c}^t = \boldsymbol{\alpha}^t Y, \text{ with } \boldsymbol{\alpha}^t = \text{softmax}(\mathbf{a}^t)$$

$$M = \text{MFB}(X, \mathbf{c}^t)$$

Visual attention

$$\mathbf{a}^{v,(r)} = \text{Convolution1D}_{v,2}^{(r)}(\text{ReLU}(\text{Convolution1D}_{v,1}^{(r)}(M)))$$

$in=d_g, out=1, k=1$ $in=2d_g, out=d_g, k=1$

$$\mathbf{c}^{v,(r)} = \boldsymbol{\alpha}^{v,(r)} M, \text{ with } \boldsymbol{\alpha}^{v,(r)} = \text{softmax}(\mathbf{a}^{v,(r)})$$

$$\mathbf{c}^v = W_f[\mathbf{c}^{v,(1)}; \mathbf{c}^{v,(2)}; \dots; \mathbf{c}^{v,(R)}]$$

Multimodal item representation

$$\text{MFB}(\mathbf{c}^v, \mathbf{c}^t)$$

12

Datasets

Polyvore68K-ND

- 68,306 outfits
(78% training, 7% validation, 15% testing)
- 365,054 items

Polyvore68K-D

- 35,140 outfits
(48% training, 9% validation, 43% testing)
- 175,485 items

Polyvore21K

- 20,925 outfits
(81% training, 6% validation, 13% testing)

	Item Types
Polyvore68K	Accessories, All body, Bags, Bottoms, Hats, Jewellery, Outerwear, Scarves, Shoes, Sunglasses, Tops
Polyvore21K	Accessories, Activewear, Baby, Bags and Wallets, Belts, Boys, Cardigans and Vests, Clothing, Costumes, Cover-ups, Dresses, Eyewear, Girls, Gloves, Hats, Hosiery and Socks, Jeans, Jewellery, Jumpsuits, Juniors, Kids, Maternity, Outerwear, Pants, Scarves, Shoes, Shorts, Skirts, Sleepwear, Suits, Sweaters and Hoodies, Swimwear, Ties, Tops, Underwear, Watches, Wedding Dresses

Table 2: Item types kept in the Polyvore68K and Polyvore21K datasets.

Experimental setup

Experiments and evaluation

- *Fashion compatibility (FC) task*: Given a set of items, compute the outfit compatibility score as the average compatibility score across all item pairs in the set



- *Fill-in-the-blank (FITB) task*: Given an incomplete set of items and 4 candidate items, find the most compatible candidate item as the one which has the highest total compatibility score with the items in the set



Training details

- Output of the 7x7x256 res4b_relu layer of ResNet18 to represent images
- Bidirectional LSTM to represent descriptions and words

Results

	Polyvore68K-ND		Polyvore68K-D		Polyvore21K	
	FC	FITB	FC	FITB	FC	FITB
<i>Common space fusion</i> baseline [11]	85.62	56.55	85.07	56.91	86.28	58.35
<i>Attention-based fusion</i>						
visual dot product attention	89.43	61.55	86.85	60.12	88.59	63.11
stacked visual attention	89.68	61.92	87.25	60.48	88.89	62.52
co-attention	89.58	61.20	86.25	59.00	85.04	58.20

Table 1: Results on the fashion compatibility and fill-in-the-blank tasks for the Polyvore68K dataset versions and the Polyvore21K dataset.

Results

FITB question:



Answers:

Baseline:



Ours:



Results

FITB question:



Answers:



Baseline:



Ours:

Results

FITB question:



Answers:



Baseline:

Ours:



Results

FITB question:

dogeared sterling silver **elephant** stud



Answers:



Baseline:



Ours:

women's j.crew for david sheldrick
wildlife trust **elephant** t-shirt

Conclusions and future work

- Attention on region-level image features and word-level text features allows to bring certain product features to the forefront in the multimodal item representations, which benefits the outfit recommendation task
- Improve state-of-the-art results on an outfit compatibility prediction task and an outfit completion task on three datasets
- Investigate neural architectures that still better recognise fine-grained fashion attributes in images
- Design novel co-attention mechanisms

