

ML FASHION ITEM RETRIEVAL

CS- 308 MACHINE LEARNING
UNDER THE GUIDANCE OF DR. JIGNESH BHATT

AMISHA LALWANI- 202251013

Citation Detail of Paper: LEARNING FASHION COMPATIBILITY ACROSS APPAREL CATEGORIES FOR OUTFIT RECOMMENDATION Luisa F. Polanía, Satyajit Gupte: [2019 IEEE International Conference on Image Processing (ICIP)]

PROBLEM STATEMENT

Goal: Assist small fashion businesses with personalized recommendation features.

Given: A query image by the user.

Objective: Learn a model that recommends compatible outfits.

Constraints: Limited training - 2000 outfits and 5000 outfits only in 80:20 train:validate ratio.

Approach/Technology Used:

- Siamese Network with VGG-16 backbone for feature extraction and color histogram fusion for better compatibility learning.
- Instead of simple concatenation used hadamard product between the image features extracted from 2 images and the two color histograms then concatenated the final results.
- Fully connected DeepFCMetricNetwork for learning the non linear compatibility relation between fashion items image.

DATASET DESCRIPTION

Dataset used: Polyvore Outfits Dataset

Train: Used train.json from non- disjoint set, split into 80:20 for training and validation. First I used only 2000 outfits to train but due to not good enough results re-trained for 5000 outfits.

Test: Used different outfits from test.json from non-disjoint set to test.

Machine Specification: Used Kaggle notebook with GPUP100.

The dataset consists of two types of sets: **disjoint** and **non-disjoint**:

In the non-disjoint set, some items (but not complete outfits) appear in both training and test splits. This allows models to leverage previously seen items to predict outfit compatibility.

In the disjoint set, the training and test splits have no overlapping items. This ensures that models must generalize to unseen outfits without relying on previously encountered fashion items.

Dataset statistics:

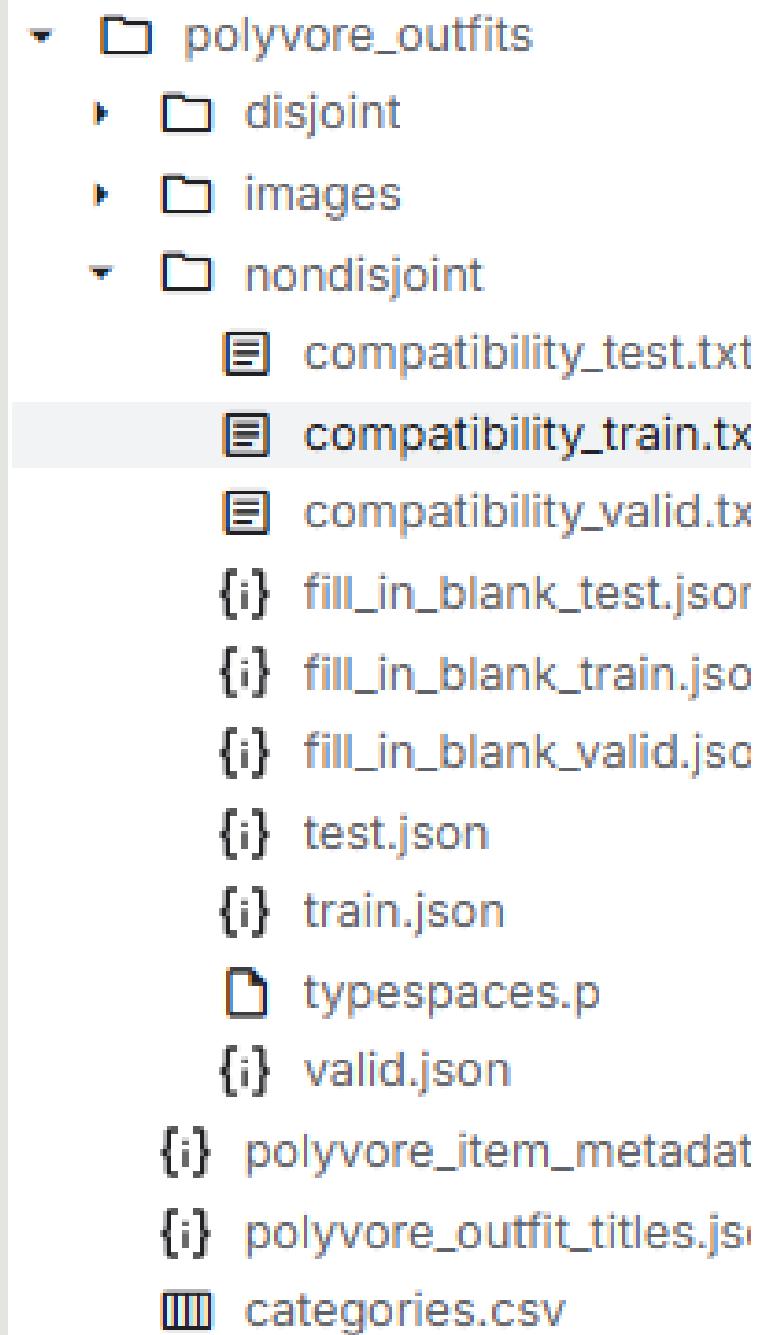
Non-disjoint set: 53,306 training, 5,000 validation and 10,000 testing outfits.

Disjoint set: 16,995 training outfits, 3,000 validation and 15,145 testing outfits.

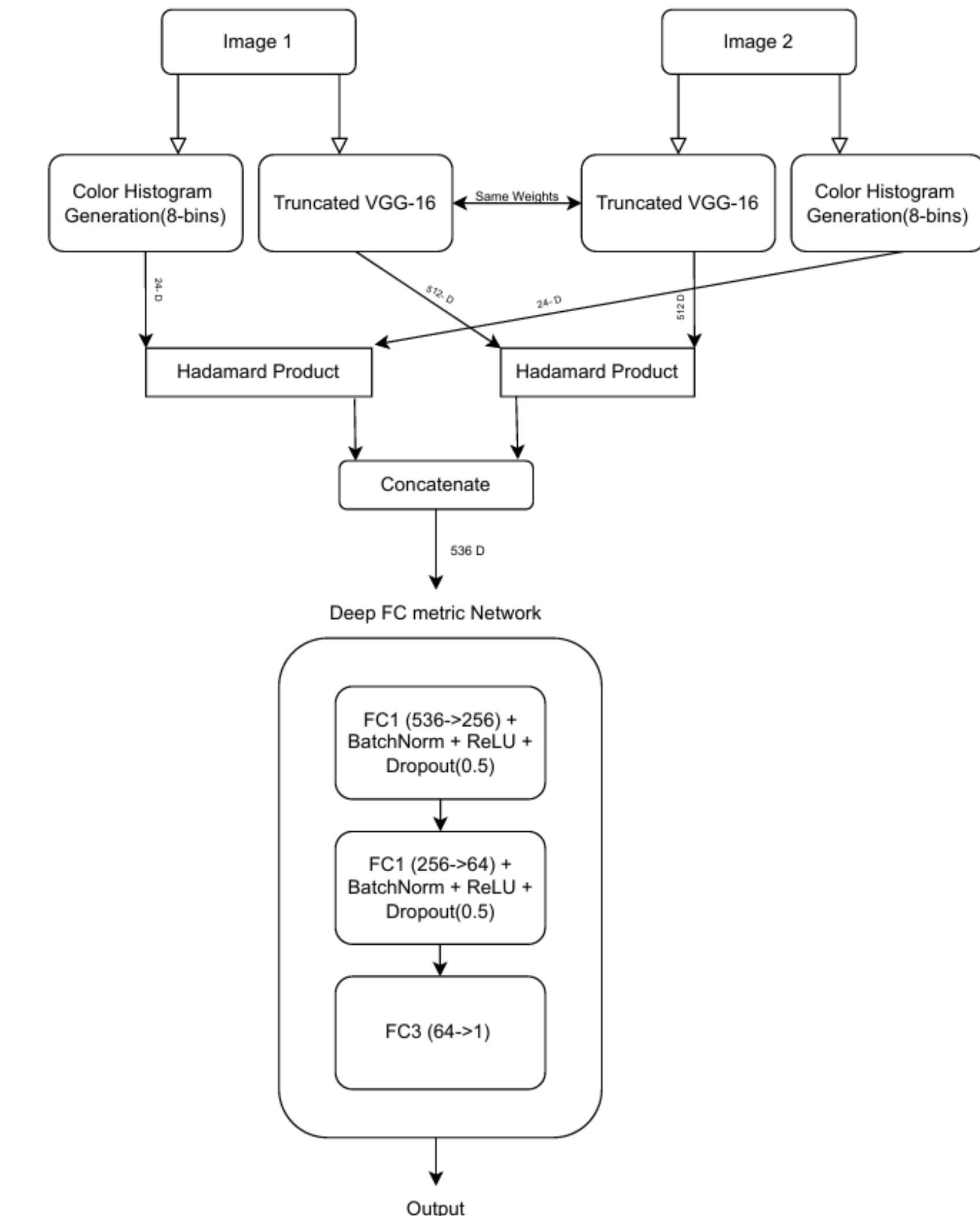
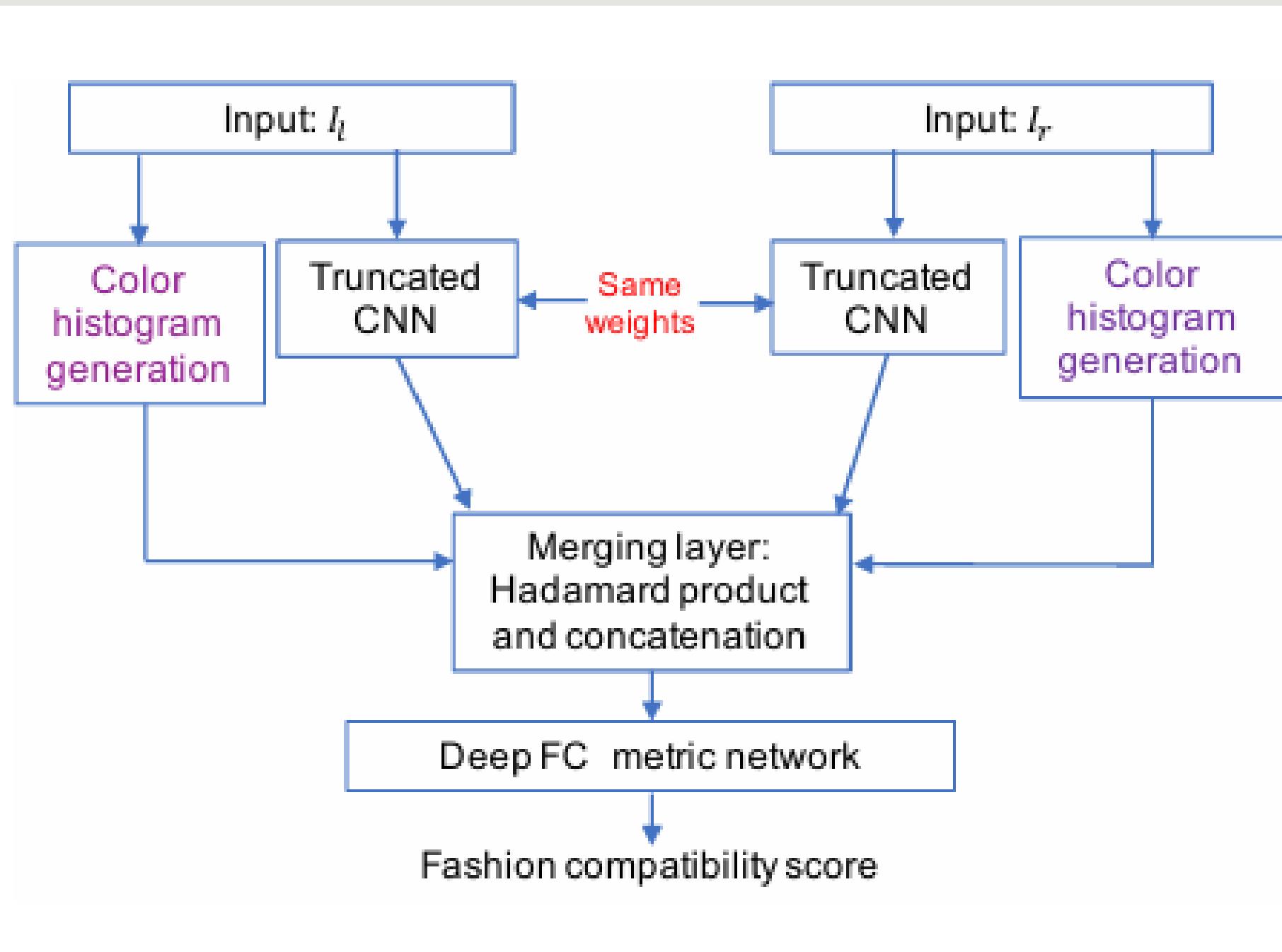
Each item has information in polyvore_item_metadata.json.

Category ids can be matched via categories.csv.

Each set has information in polyvore_outfit_titles.json.

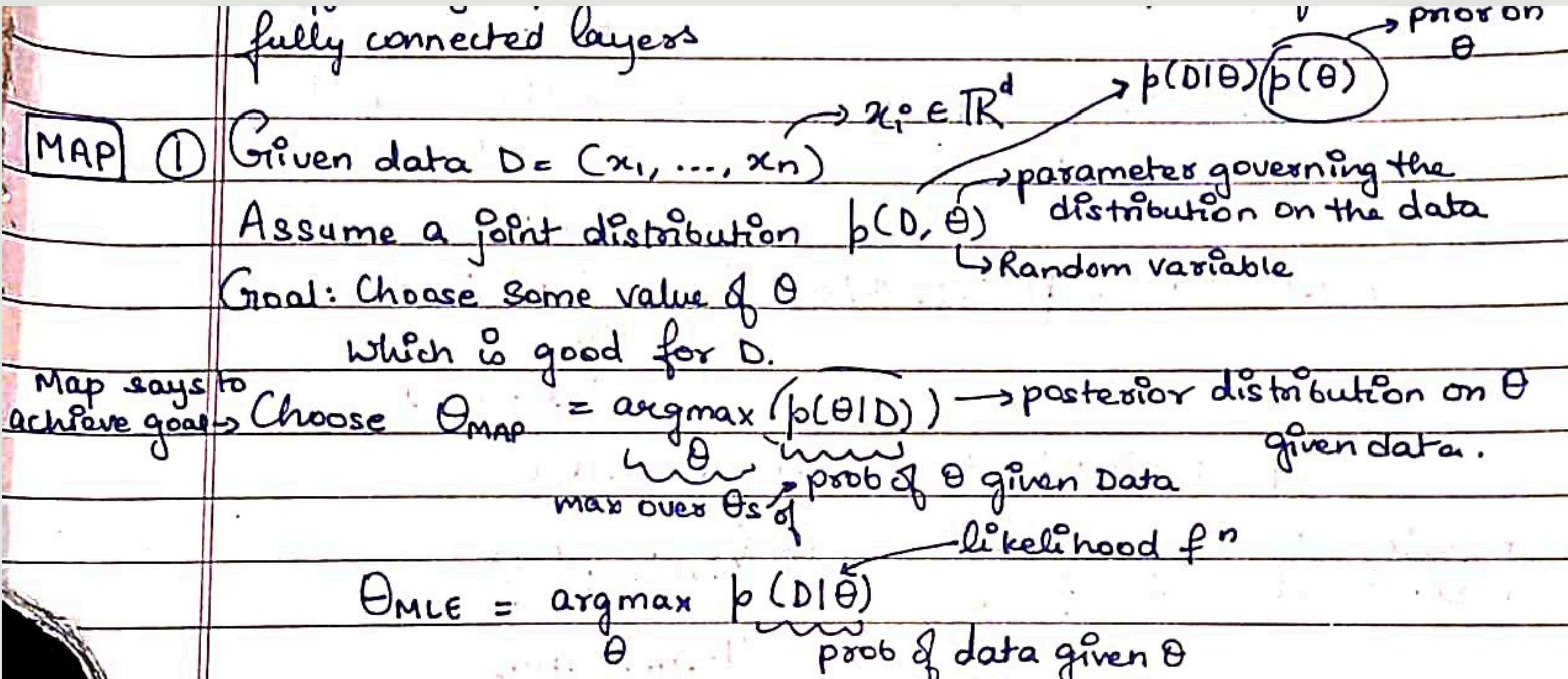


NEURAL NETWORK ARCHITECTURE



KEY TERMS TO KNOW

MAP - MAXIMUM APOSTERIORI PROBLEM



④ Hadamard product: Element-wise product

$$A = \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix} \quad B = \begin{bmatrix} 7 & 8 \\ 9 & 10 \end{bmatrix} \quad A \odot B = \begin{bmatrix} 1 & 2 \\ 4 & 5 \end{bmatrix} \odot \begin{bmatrix} 7 & 8 \\ 9 & 10 \end{bmatrix} = \begin{bmatrix} 7 & 16 \\ 36 & 50 \end{bmatrix}$$

- Captures interaction btw corresponding dimensions. (like matching features)
- It helps to see what features overlap the most for compatibility

TRAINING PROCESS

N training I/P image pairs $I = \{(I_r, I_e)_i\}_{i=1}^N$

I_e : I/P to left] branches of siamese sub-network.

I_r : I/P to right] $y = \{y_i\}_{i=1}^N$ denotes binary labels

$y_i = 1$ if I/P pair $(I_r, I_e) \rightarrow$ fashion compatible
else 0.

x_i - o/p of last fc layer of metric sub-network

when pair (I_r, I_e) - i/p. Readout fn $\gamma(\cdot)$ is

$$\hat{y}_i = \gamma(x_i) = w^T x_i + \epsilon \rightarrow \text{Noise has logistic distribution}$$

wts. of
readout fn

standard logistic
distribution.

Batch Normalization applied after every FC layer of metric sub networks

Siamese sub net 2 set of wts Θ_L, Θ_R - but mirrored $\therefore \Theta = \Theta_L = \Theta_R$

$\Theta_S = \{\Theta_t\}_{t=1}^S \rightarrow$ subset of Θ that correspond to S filters selected for fine tuning.

$W = \{w_j \in \mathbb{R}^{P_j \times Q_j}\}_{j=1}^M$ be the set of wts of fc metric sub networks.

→ Modeled with a Matrix-variate normal distribution of 0 mean, i.e. $w_j \sim M.N(0, \Lambda_j, \Sigma_j^2 I)$ $\forall j \quad 0 \in \mathbb{R}^{P_j \times Q_j}$

0 matrix &

$\Sigma_j^2 I \in \mathbb{R}^{P_j \times P_j}$

row covariance diag matrix.

The positive semi-definite col cov matrix $\in \Lambda_j \in \mathbb{R}^{Q_j \times Q_j}$

covariance matrix - can be learned in order to

Capture correlation btw layer i/p units

$\Lambda = \{\Lambda_j\}_{j=1}^M$ - set of col cov matrix

Network is trained by solving the following MAP problem

$$\hat{\theta}_s, \hat{W}, \hat{\Lambda}, \hat{w} = \underset{\theta_s, W, \Lambda, w}{\operatorname{argmax}} p(\theta_s, W, \Lambda, w | I, Y)$$

$$= \underset{\theta_s, W, \Lambda, w}{\operatorname{argmax}} p(Y | I, \theta_s, W, w) \times \textcircled{3}$$

likelihood $\xleftarrow{\text{prior probabilities}}$ $p(W | \Lambda) \times p(\theta_s) \times p(w)$
 probability $\xrightarrow{\text{prior probabilities}}$ $p(\theta_s) \times p(w)$

noise (ϵ) - logistic
 prob distrib' of y_i given x_i - Bernoulli \therefore likelihood

$$p(Y | I, \theta_s, W, w) \propto \prod_{i=1}^N p(y_i = 1 | x_i) \overset{y_i}{p}(1 - \overset{1-y_i}{p}(y_i = 1 | x_i)) \textcircled{4}$$

$$p(y_i = 1 | x_i) \overset{y_i}{p} = \sigma(w^T x_i) =$$

$$\frac{1}{1 + \exp\{-w^T x_i\}}$$

Each entry of the vectors in the set Θ & each entry of w is modeled with Laplacian distribution of zero mean & variance σ_p^2 , in case of Θ & variance σ_w^2 in case of w . Motivation is to promote sparsity i.e. reduce redundancy in the wts.



Selected CNN filters

Laplacian($0, \sigma_p^2$) Promotes sparse filters

Readout weights

Laplacian($0, \sigma_w^2$) Promotes sparse decision logic

w_i - modeled by matrix-variate normal distribution with 0 mean, the prior prob $p(w|\Lambda)$

takes the form

M

$$p(w|\Lambda) = \prod_{j=1}^M p(w_j|\Lambda_j) \quad (5)$$

tr(.) - trace of matrix

|.| - determinant

$$\frac{\exp(-\frac{1}{2} \text{tr}((\Sigma_j^{-1})^{-1} w_j \Lambda_j^{-1} w_j^T))}{(2\pi)^{\frac{p_j Q_j}{2}} |\Sigma_j^{-1}|^{\frac{p_j}{2}} |\Lambda_j|^{-\frac{Q_j}{2}}} \quad (6)$$

Replacing 4,6 in 3 and taking the negative logarithm leads to the following optimization problem

$$-\text{ve log} : \arg\max p(\text{likelihood}) \cdot p(\text{prior}) \Rightarrow \arg\min -\log(p(\cdot))$$

$$\underbrace{\sum_{i=1}^N [-y_i^\circ \log(w^\top x_i) - (1-y_i^\circ) \log(1 - \tau(w^\top x_i^\circ))]}_{\text{-ve log likelihood}}$$

$$\underbrace{\text{Cross entropy loss from eqn 4}}_{\text{Likelihood}}$$

$$+ \sum_{j=1}^M \frac{1}{2} \text{tr}((\lambda_j^{-2} I)^{-1} w_j^\circ \lambda_j^{-1} w_j^{0T}) + \frac{\rho_j^\circ \log |\lambda_j|}{2}$$

from Eqn 6.

$$+ \frac{1}{2} \|\Theta_s\|_1 + \frac{1}{2} \|\boldsymbol{\omega}\|_1$$

Laplacian prior.

LOSS FUNCTION

Term-1 : BCE LOSS-

Helps the model distinguish between compatible and incompatible pairs.

Term-2: Matric-Normal-Prior-

Regularizes Weight matric W penalizes high variance in weights promotes structured regularization guided by Λ_j - encourages weights to learn coherent embeddings

Term-3: CNN Sparsity Prior

Promotes sparsity in CNN filter weights, L1 norm makes many small weights to be 0 → prevents overfitting, improves generalization smaller inference time.

$$\hat{\Theta}_s, \hat{\mathbf{W}}, \hat{\Lambda}, \hat{w} = \underset{\Theta_s, \mathbf{W}, \Lambda, w}{\arg \min} - \sum_{i=1}^N [y_i \ln(r(w^T x_i)) + (1 - y_i) \ln(1 - r(w^T x_i))] + \lambda_1 \sum_{j=1}^M \text{tr}(W_j \Lambda_j^{-1} W_j^T) + \lambda_2 \sum_{t=1}^S \|\theta_t\|_1 + \lambda_3 \|w\|_1 \text{ s. t. } \Lambda_j \succeq 0, \text{tr}[\Lambda_j] = 1, \forall j,$$

where $\lambda_1, \lambda_2, \lambda_3$ are regularization parameters, which are incorporated to tune the strength of the regularization terms, and are estimated using grid search. The constraint $\Lambda_j \succeq 0$ comes from the positive semi-definite property that covariance matrices need to satisfy. An alternating optimization procedure

Term-4: Classifier Sparsity Prior

Forces feature selection by suppressing irrelevant features-> makes decision boundary rely only on key features.

RESULTS- QUALITATIVE AND QUANTITATIVE

To measure result quantitatively I used Precision@K and Lift@K as explained in the paper.

Let $\{\psi_n\}_{n=1}^{N_t}$ be the set of testing Polyvore outfits, where ψ_n is formed by the query item (first item in the outfit), denoted as q_n , and C_n complementary items, denoted as $\{o_n^{(c)}\}_{c=1}^{C_n}$, which belong to C_n different apparel categories. Let $R_n^c(K)$ denote the top K recommendations generated by the proposed network for the complementary item in category c , given the query q_n . To generate $R_n^c(K)$, pairs are first formed between the query item and all the rest of the items in the Polyvore dataset which belong to category c , then the items with the top K fashion compatibility scores are selected to form $R_n^c(K)$. The precision@ K for outfit ψ_n is

$$\text{precision}@K(\psi_n) = \frac{1}{C_n} \sum_{c=1}^{C_n} \mathbf{1}[o_n^{(c)} \in R_n^c(K)], \quad (10)$$

where $\mathbf{1}[\cdot]$ denotes the indicator function. The average of the precision@ K across the N_t testing outfits is referred to as the average precision@ K . The recommendation performance of a model is evaluated using the lift of average precision@ K , which is defined as

$$\text{Lift}@K = \frac{\text{average precision}@K(\text{model})}{\text{average precision}@K(\text{random})}, \quad (11)$$

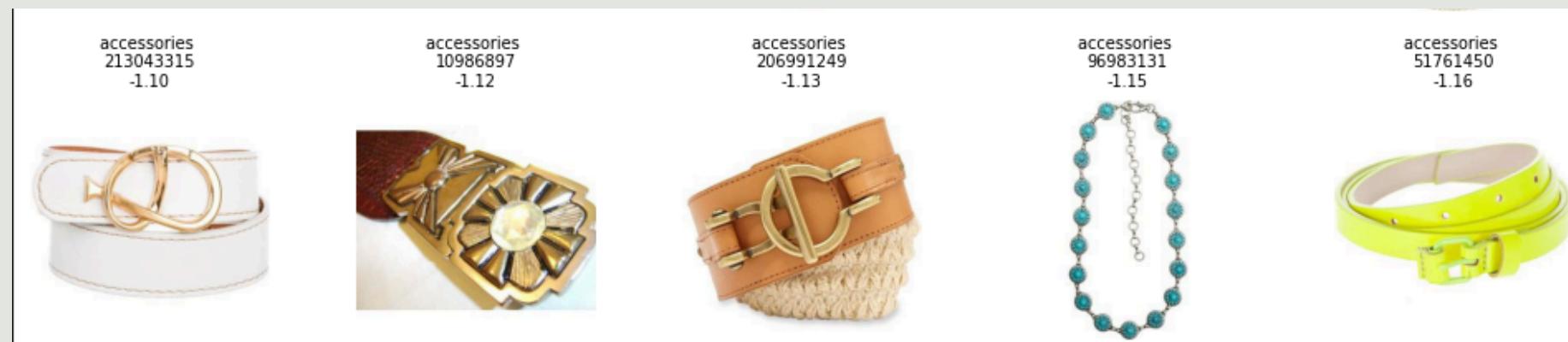
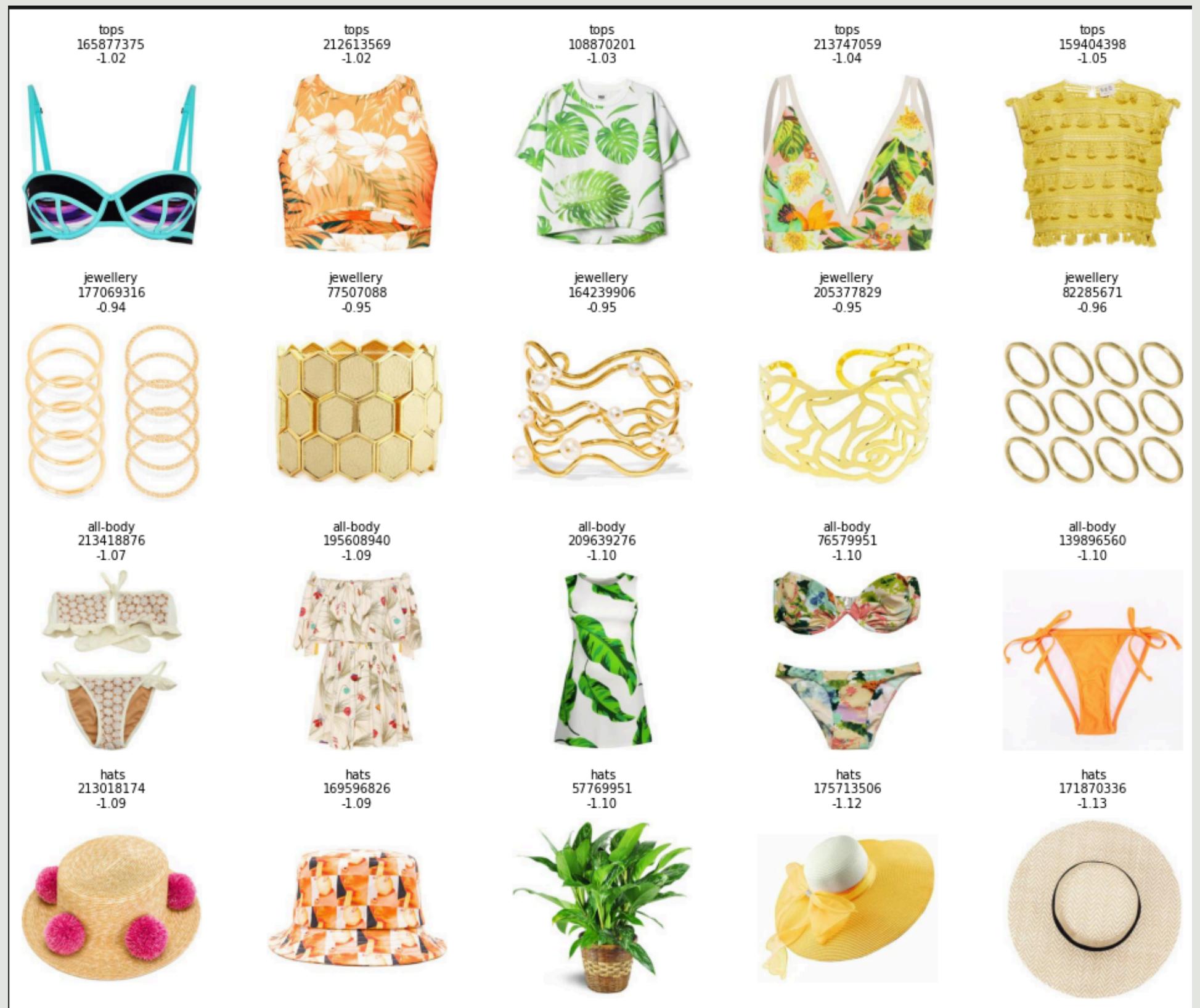
where average precision@ K (random) is that of a recommender that would select items at random for each of the apparel categories of interest.

Precision@K just represents how many outfits from the real outfit set of query get recommended by my model.

RESULTS- QUALITATIVE AND QUANTITATIVE

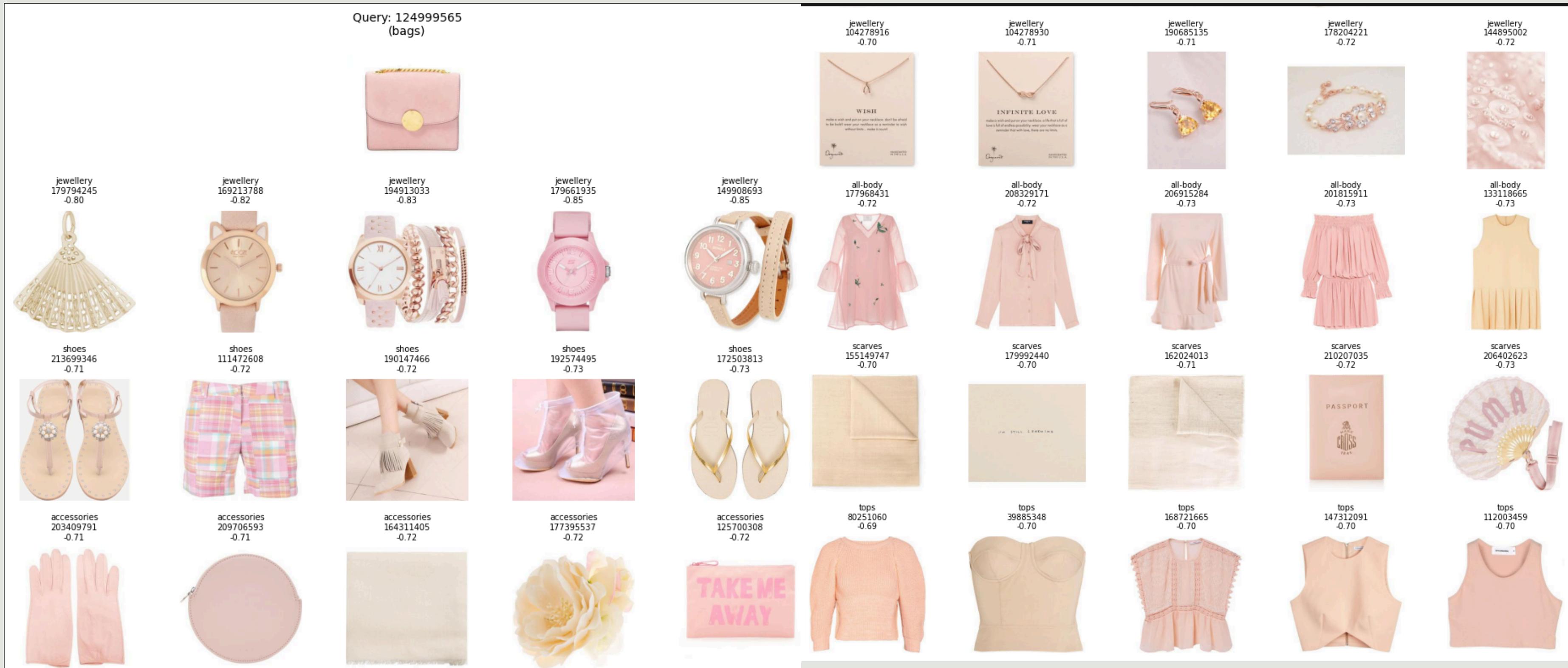
When model was trained for 1600 outfits from train.json and validated for 400 outfits from train.json. One sample result is shown:-

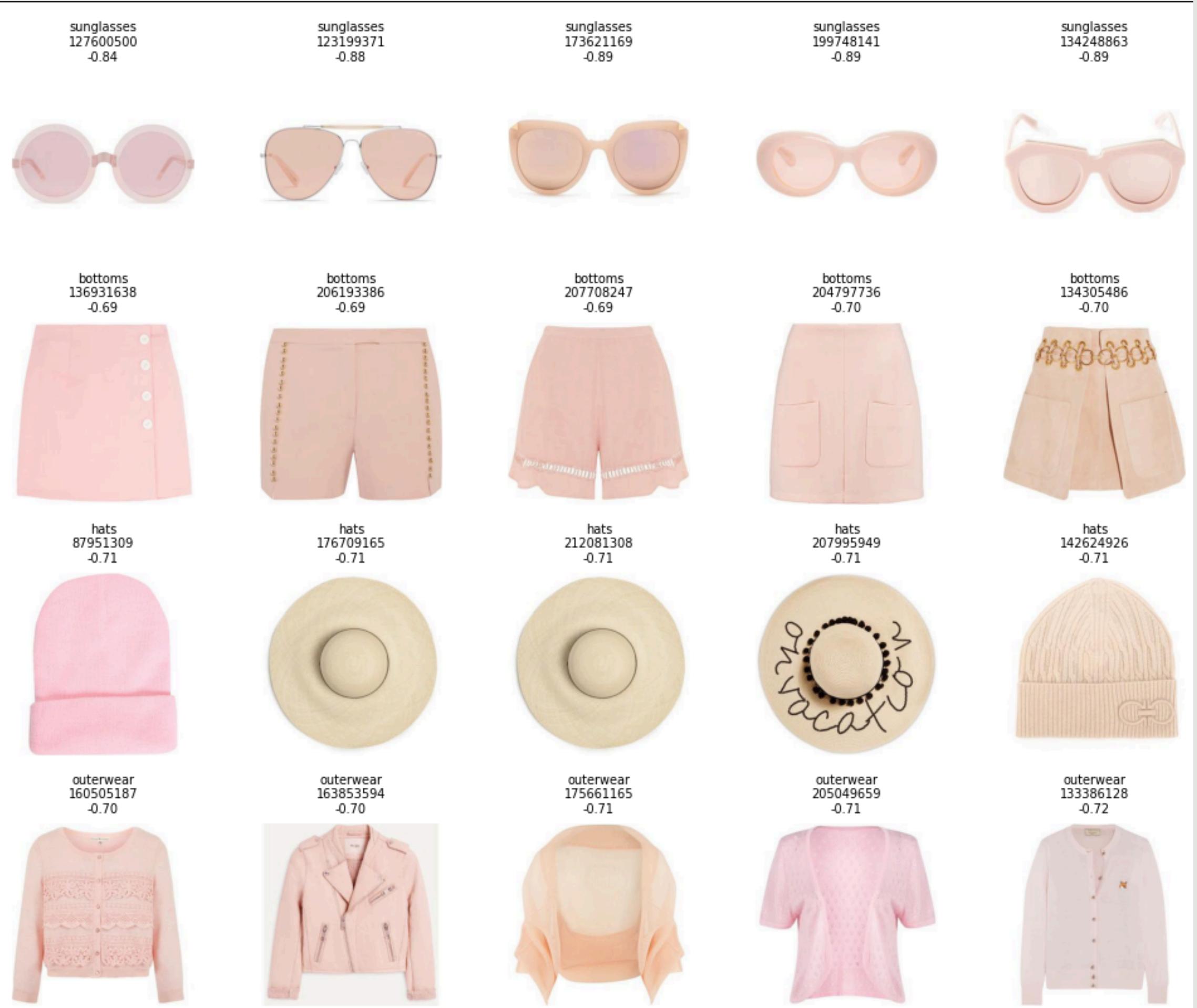




Precision@5 came out to be zero i.e. none of the recommended items belonged to the original outfit set of the query, in paper they used 10650, 1902, 1395 for training, validation and testing so I had used only 1600 outfits for training and 400 for validation and 5 for testing so I increased training to 4000 outfits validation to 1000 outfits.

When model was trained for 4000 outfits from train.json and validated for 1000 outfits from train.json. One result for same query is shown we can see the improvements in the recommended items :-





But Precision@K
still was 0. So I
increased K=5 to
K=100 then also
Precision@K= 0

Cause of low precision

This low precision@K might be due to the fact that I am not taking Lambda as defined in the paper, I used $\Lambda_j = \|W_j\|^2$ then in the loss function trace term becomes=1 constant meaning all input features are equally important and uncorrelated while acc to the closed form definition of lambda it is that some features are correlated- shrink the directions that are not relevant. Also the training set is also small I need to change update_lambda defination as well as train on more data.

Thank You

For your attention