

# Deep Text Mining of Instagram Data Without Strong Supervision

Master's Thesis Defense

Kim Hammar

KTH Royal Institute of Technology

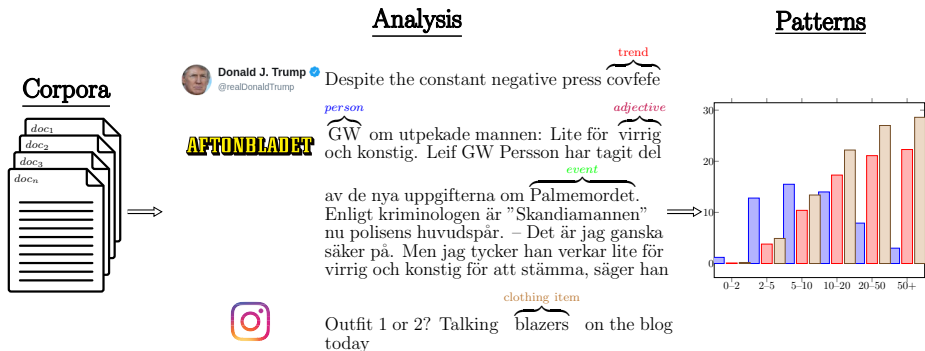
*kimham@kth.se*

June 1, 2018



# What Is Text Mining?

Extracting information and detecting patterns in unstructured text



Text understanding is hard → AI-Complete Problem

- **Ambiguous** sentences: we need context to understand

*Did you see her dress?*

- **Ambiguous** sentences: we need context to understand

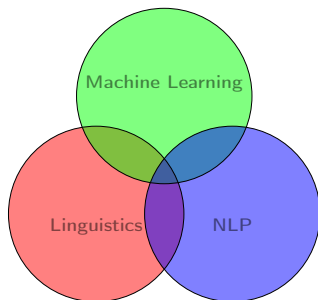
*Did you see her dress?*

**Yes I was in the hall with her when she dressed.**

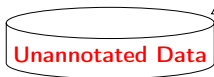
**Yes It was gorgeous!**

- Classic NLP & Text Mining: **Linguistic Rules**
  - “Two capital words in a row” → A person’s name
- Data Driven Text Mining: **Machine Learning Models**

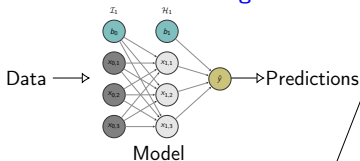
$$\nabla_{\theta} L = \left( \frac{\partial L}{\partial W_{1,1}}, \dots, \frac{\partial L}{\partial W_{n,n}} \right)$$



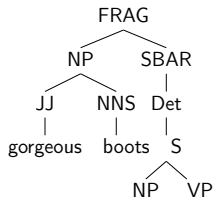
# Unsupervised Text Mining in Social Media



## Machine Learning



## Natural Language Processing (NLP)



## Domain Under Study: Fashion



Automatic predictions:

Items: jeansshorts, ...

fabrics: denim, ...

style: casual, ..

⋮



“Noisy” Text  
outfit of the  
day 🥰 #goals

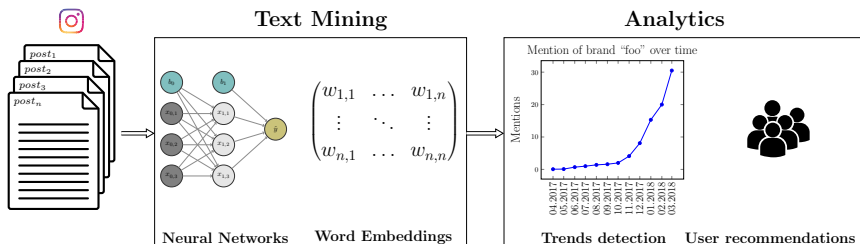
## Big Data

40B images

500M daily activity

# Problem: Unsupervised Text Mining from Instagram

- Input: Noisy Text (Image Captions, User Comments)
- Output: Fashion Attributes (Items, Fabrics, Brands... etc.)



# Example Instagram Post



[Redacted]

Following

[Redacted] Wow 🤩

[Redacted] C'est le même que toi 😊

[Redacted] Where is that coat from??! I have to have it!!

[Redacted] It's soooo fluffy I'm gonna diiiiie  
😊💕

[Redacted] So cozy and noticeable 🙌💕

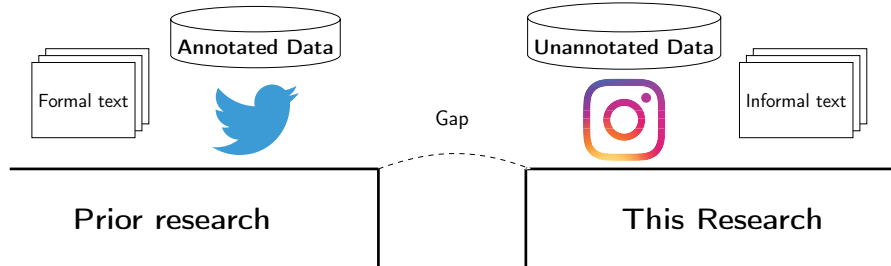
[Redacted] [@hugoboss](#) [@hugoboss](#) H&M! Sold out tho 😞

[Redacted] Beauty 🙌🙌🙌🙌

[Redacted] Bästa jackan!! 😊 Så glad jag fick tag i en också för ett tag sen 😊 så fruktansvärt varm dock, måste nästan vänta på snö... [Redacted]



## How Our Research Stands Out From Prior Work



- 1 **An empirical study of Instagram text**
  - No previous study on Instagram text exists that we are aware of
- 2 **Unsupervised extraction of fashion attributes from Instagram using Word Embeddings**
  - The first evaluation of word embeddings for Instagram
  - The first distributed implementation of the FastText algorithm
  - We confirm prior results on IE and apply it to a new domain
- 3 **Novel pipeline for classification with weak supervision & deep learning**
  - Extension of the data programming paradigm to the multi-label setting

---

<sup>0</sup>All code and most of the data is open source:

<https://github.com/shatha2014/FashionRec>

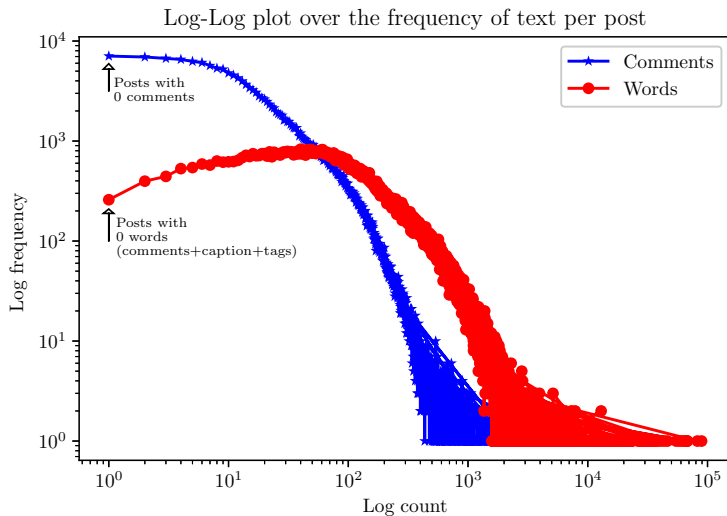
# Instagram Text is Noisy and Multi-Lingual

A case study of a corpora with 143 fashion accounts, 200K posts, 9M comments

- **Instagram text is noisy:** 47% OOV words when including URLs, emojis etc. Otherwise 30% (compared to 25% on Twitter)
- **Comment sections are multi-lingual:** All accounts are English, still only 52% of comments are English (total 97 languages identified)
- **The text is ungrammatical:** Informal spelling, unreliable capitalization.

# Instagram Text Distribution Has a Long Tail

A case study of a corpora with 143 fashion accounts, 200K posts, 9M comments

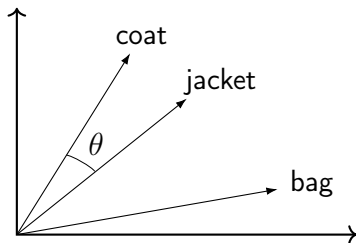


# Conclusion from Case Study: Prefer Statistical Methods Rather Than Symbolic NLP Methods

- **Instagram Text is noisy, multi-lingual, and un-grammatical**
  - → Linguistic methods for text mining are fragile
  - → Syntactic text matching is difficult (many languages, many synonyms, online-specific tokens etc)
- **We propose: Word Embeddings as a key component in information extraction from Social Media**
  - Demonstrated in the second contribution of the thesis

# Word Embeddings Are Distributed Representation of Words

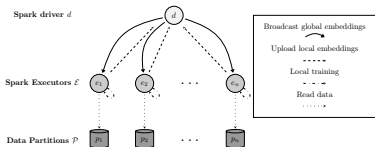
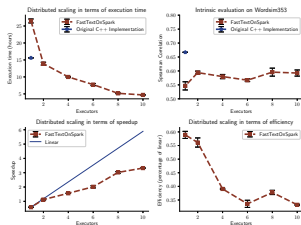
- Word Embeddings are vectors in  $\mathbb{R}^d$ ,  $d \approx 300$
- Derived with optimization using the *Distributional Hypothesis*<sup>1</sup>
  - → Words that occur in similar contexts will obtain similar vectors
  - “You shall know a word by the company it keeps” - Firth '57<sup>2</sup>
  - → we can use the word knowledge in word embeddings for IE



<sup>1</sup>Zellig S Harris. “Distributional structure”. In: *Word* 10.2-3 (1954), pp. 146–162.

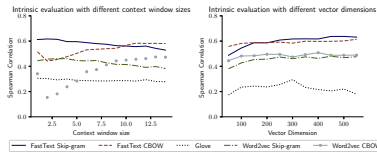
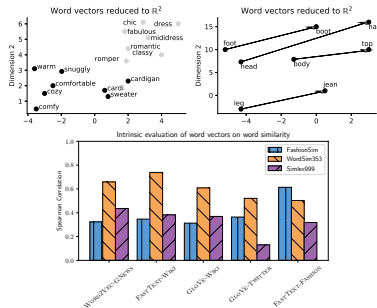
<sup>2</sup>J. R. Firth. “A synopsis of linguistic theory 1930-55.” In: 1952-59 (1957), pp. 1-32. <img alt="Navigation icons: back, forward, search, refresh, etc."/>

# Contributions On Word Embeddings



## FastTextOnSpark:

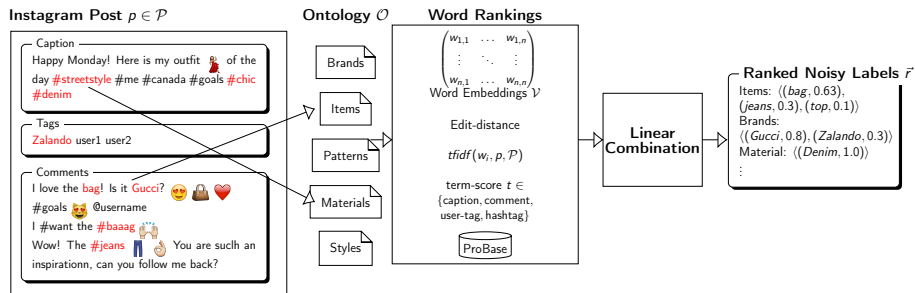
- A Scalable Implementation of FastText



## Instagram Word Embeddings:

- Hyperparameter Tuning
- Comparison With Off-The-Shelf Embeddings
- Extrinsic and Intrinsic Evaluations

# Unsupervised Information Extraction using Word Embeddings and a Fashion Ontology



**Figure:** An information extraction system for social media text. The system extracts fashion details from text associated with Instagram posts.



# Word Embeddings Outperform Syntactic Baseline for IE

## Results

- IE with word embeddings **outperform** IE based on edit-distance ( $p < 0.05$ )

<i>Method/Category</i>	<i>NDGC@1</i>	<i>NDGC@3</i>	<i>NDGC@5</i>	<i>NDGC@10</i>	<i>P@1</i>	<i>P@3</i>	<i>P@5</i>	<i>P@10</i>	<i>MAP</i>
SemCluster/Item	<b>0.833</b>	<b>0.658</b>	<b>0.691</b>	<b>0.807</b>	<b>0.833</b>	<b>0.546</b>	<b>0.454</b>	<b>0.309</b>	<b>0.733</b>
SynCluster/Item	0.781	0.581 <sup>-</sup>	0.607 <sup>-</sup>	0.767 <sup>-</sup>	0.781	0.474 <sup>-</sup>	0.370 <sup>-</sup>	0.296	0.641 <sup>-</sup>
SemCluster/Style	<b>0.399</b>	<b>0.505</b>	<b>0.519</b>	<b>0.548</b>	<b>0.417</b>	<b>0.204</b>	<b>0.139</b>	<b>0.069</b>	<b>0.539</b>
SynCluster/Style	0.367	0.415 <sup>-</sup>	0.425 <sup>-</sup>	0.507	0.367	0.130 <sup>-</sup>	0.123	0.069	0.474 <sup>-</sup>
SemCluster/Pattern	0.087	0.179	0.353	0.444	0.087	0.110	0.169	<b>0.118</b>	0.296
SynCluster/Pattern	<b>0.108</b>	<b>0.413</b>	<b>0.498</b>	<b>0.512</b>	<b>0.108</b>	<b>0.221</b>	<b>0.193</b>	<b>0.117</b>	<b>0.395</b>
SemCluster/Material	<b>0.296</b>	<b>0.286</b>	<b>0.324</b>	<b>0.393</b>	<b>0.286</b>	<b>0.264</b>	<b>0.233</b>	<b>0.165</b>	<b>0.373</b>
SynCluster/Material	0.113 <sup>-</sup>	0.104 <sup>-</sup>	0.137 <sup>-</sup>	0.209 <sup>-</sup>	0.113 <sup>-</sup>	0.107 <sup>-</sup>	0.109 <sup>-</sup>	0.092 <sup>-</sup>	0.227 <sup>-</sup>
SemCluster/Brand	<b>0.062</b>	<b>0.066</b>	<b>0.062</b>	<b>0.064</b>	<b>0.032</b>	<b>0.056</b>	<b>0.036</b>	<b>0.039</b>	<b>0.194</b>
SynCluster/Brand	0.016	0.010	0.010	0.010	0.016	0.005	0.003	0.002	0.159

<sup>2</sup>A smaller, hand-labeled dataset by experts was used for evaluation

## Text Mining with Word Embeddings and a smaller domain ontology:

### Pros:

- Does not require annotated data
- Can deal with noisy text
- Transparent model

### Cons:

- Require feature engineering
- Require a domain ontology

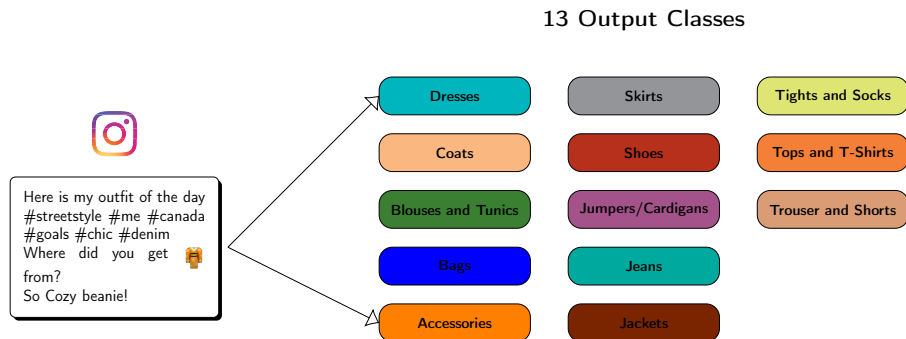
*Can we reduce manual feature engineering and learn from data?*

Problem: **We don't have annotated data (yet)**

Solution: **Weakly supervised learning**

# The Classification Task

## A multi-class multi-label classification problem



# Alternative Sources of Supervision That Are Cheap but Weak

- **Strong supervision:** Manual annotation by expert
- **Weak supervision:** A signal that does not have full coverage/perfect accuracy

## Sources of Weak Supervision



Domain Heuristics



Crowdworkers



# How To Combine Several Sources Of Weak Supervision?

- Simplest way to combine many weak signals: **Majority Vote**
- Recent research on combination of weak signals: **Data Programming Paradigm**<sup>3</sup>


---

## Data Programming: Creating Large Training Sets, Quickly

---

Alexander Ratner, Christopher De Sa, Sen Wu, Daniel Selsam, Christopher Ré  
Stanford University  
{ajratner, cdesa, senwu, dselsam, chrismre}@stanford.edu

---

<sup>3</sup>Alexander J Ratner et al. "Data Programming: Creating Large Training Sets, Quickly". In: *Advances in Neural Information Processing Systems 29*. Ed. by D. D. Lee et al. Curran Associates, Inc., 2016, pp. 3567–3575. url: <http://papers.nips.cc/paper/6523-data-programming-creating-large-training-sets-quickly.pdf>. 

# Data Programming: Model Weak Supervision With Generative Model

- **Model weak supervision as labeling functions**  $\lambda_i$ 
  - $\lambda_i(\text{unlabeled data}) \rightarrow \text{label}$
- **Learn Generative Model**  $\pi_{\alpha, \beta}(\Lambda, Y)$  **over the labeling process.**
  - Based on conflicts between **labeling functions** assign the functions an **estimated accuracy**  $\alpha_i$ .
  - Based on empirical coverage of **labeling functions** assign the functions a **coverage**  $\beta_i$ .
- **Given  $\alpha$  and  $\beta$  for each labeling function, it can be used to combine labels into a single probabilistic label**
  - Give more weight to high-accuracy functions
  - If there is a lot of disagreement  $\rightarrow$  **low probability label**
  - If all labeling functions agree  $\rightarrow$  **high probability label**

## Low accuracy labeling functions

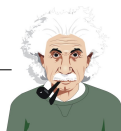


"it is not a coat"

## High accuracy labeling functions



"it is a coat"



**Probabilistic Label:** 0.6 probability that it is a coat

**Majority Vote:** 1.0 probability that it is not a coat

# Pipeline for Weakly Supervised Classification in Instagram

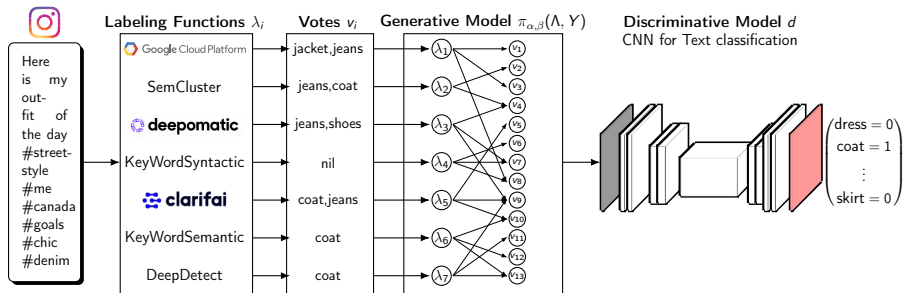


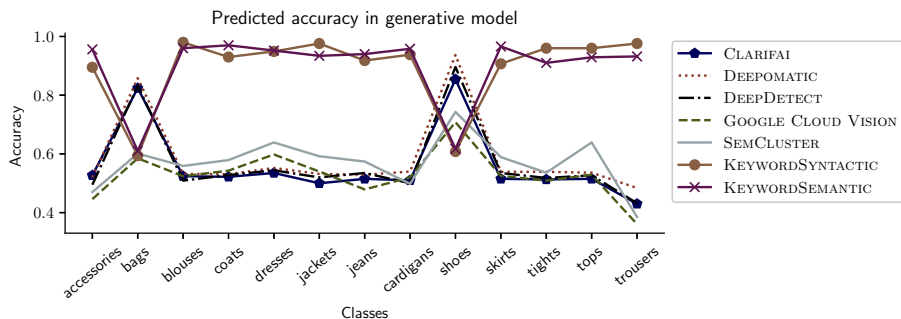
Figure: A pipeline for weakly supervised text classification of Instagram posts.



# Extension of Data Programming to Multi-Label Classification

- **Problem:** Data programming only defined for binary classification in original paper
- **To make it work for multi-class setting:** model labeling function as  $\lambda_i \rightarrow k_i \in \{0, \dots, N\}$  instead of  $\lambda_i \rightarrow k_i \in \{-1, 0, 1\}$ .
- **Idea 1 for multi-label:** model labeling function as  $\lambda_i \rightarrow \vec{k}_i = \{v_0, \dots, v_n\} \wedge v_j \in \{-1, 0, 1\}$
- **Idea 2 for multi-label:** learn a separate generative model for each class, and let each labeling function give binary output for each class  $\lambda_{i,j} \rightarrow k_{i,j} \in \{-1, 0, 1\}$ .

# Trained Generative Models: Labeling Functions' Accuracy Differ Between Classes



**Figure:** Multiple generative models can capture a different accuracy for labeling functions for different classes.


# Discriminative CNN Model for Text Classification

- I have extended Kim Yoon's CNN model for text classification<sup>4</sup>
- To train the model with probabilistic labels produced by generative model, I use a *noise-aware loss function*<sup>5</sup>:

$$\frac{1}{N} \sum_{i=0}^N -(p(Y_i|\Lambda_i) \log(\sigma(\hat{y}_i)) + ((1 - p(Y_i|\Lambda_i)) \log(1 - \sigma(\hat{y}_i)))) \quad (1)$$

---

<sup>4</sup>Yoon Kim. "Convolutional Neural Networks for Sentence Classification". In: *EMNLP. ACL, 2014*, pp. 1746–1751.

<sup>5</sup> $N$  is the number of classes,  $p(Y_i|\Lambda_i)$  is the probabilistic labels for class  $i$ , and  $\hat{y}_i$  is the logits for class  $i$  

# Data Programming Beats Majority Voting and the Multi-Channel Model Was Not Useful

## Results

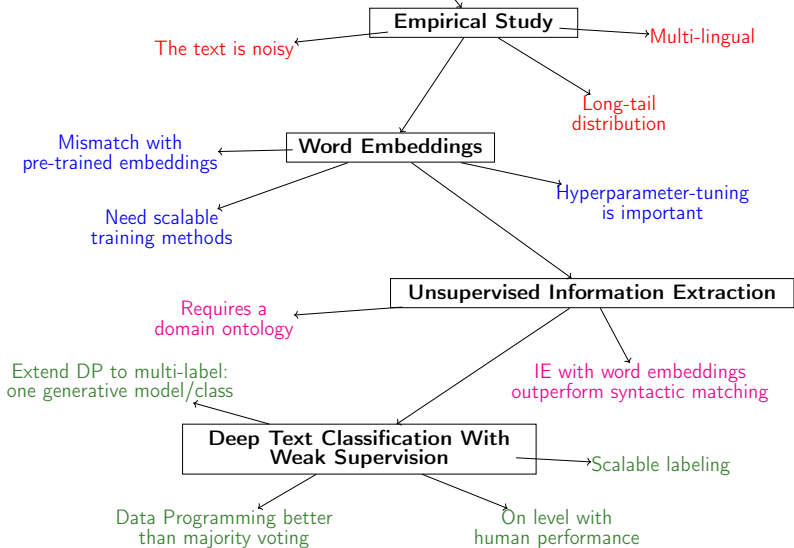
- **Data programming** gives 6  $F_1$  points improvement over majority vote<sup>6</sup>, achieving an  $F_1$  score of 0.61 (**On level with human performance**)

<i>Model</i>	<i>Accuracy</i>	<i>Precision</i>	<i>Recall</i>	<i>Micro-<math>F_1</math></i>	<i>Macro-<math>F_1</math></i>	<i>Hamming Loss</i>
CNN-DataProgramming	<b>0.797</b> $\pm$ 0.01	<b>0.566</b> $\pm$ 0.05	0.678 $\pm$ 0.04	<b>0.616</b> $\pm$ 0.02	<b>0.535</b> $\pm$ 0.01	<b>0.195</b> $\pm$ 0.02
CNN-MajorityVote	0.739 $\pm$ 0.02	0.470 $\pm$ 0.06	<b>0.686</b> $\pm$ 0.05	0.555 $\pm$ 0.03	0.465 $\pm$ 0.05	0.261 $\pm$ 0.03

- Main cause of error: **data sparsity** (can not extract clothing items from the text if it is never mentioned in the text)

<sup>6</sup>A smaller, hand-labeled dataset by experts was used for evaluation

# How Can We Do Unsupervised Text Mining From Instagram?



- Instagram text is **just as noisy as Twitter**, comment sections are multi-lingual, long tail text distribution
- **Word Embeddings** are useful for IE, especially in social media
- **Deep learning with weak supervision and data programming** is a promising approach for text mining in social media

## Main line of future work:

Combine text analytics with image analysis<sup>7</sup>

## Thanks:

Shatha Jaradat, Nima Dokoohaki Ph.D, Prof. Mihhail Matskin

---

<sup>7</sup>Shatha Jaradat. "Deep Cross-Domain Fashion Recommendation". In: *Proceedings of the Eleventh ACM Conference on Recommender Systems. RecSys '17. Como, Italy: ACM, 2017, pp. 407–410. isbn: 978-1-4503-4652-8. doi: 10.1145/3109859.3109861. url: <http://doi.acm.org/10.1145/3109859.3109861>.*

# Questions