

## The Frontier of Autoencoder-Based Recommender Systems

Martin Spišák, Vojtěch Vančura 11 December 2024



#### We Move Quickly and Adjust to the Market Needs

- 10+ years of continuous research and development
- 8 years on the market
- Exponential revenue growth







## **Outline**



- 1. How autoencoders work in collaborative filtering tasks
- 2. How to scale them to huge datasets
- 3. How to incorporate **item information**
- 4. How we **teach autoencoders to generalize** to new items



#### 1. How autoencoders work in collaborative filtering tasks

## **Interaction Autoencoders**



- Autoencoder = self-supervised representation learning technique
- In RecSys context: Train a ML model to reconstruct rows of (sparse) interaction matrix X
- Denoising AutoEncoders (DAE), Variational AutoEncoder (VAE) etc.



## EASE



- EASE is a **linear autoencoder** model with closed-form solution
  - linear regression but with huge model capacity
  - Encoder and decoder fused together
- EASE trains item-to-item weight matrix
- Diagonal of weights constrained to zero to prevent trivial solutions



## **EASE - results**



- Closed-form solution => very fast training
- (more-or-less) SotA recommendation quality

Table 1: The closed-form dense solution  $\hat{\mathbf{B}}^{(\text{dense})}$  (see Eq. 4) obtains competitive ranking-accuracy while requiring only a small fraction of the training time, compared to the various models empirically evaluated in [33]. The standard errors of the ranking-metrics are about 0.002, 0.001, and 0.001 on *ML-20M*, *Netflix*, and *MSD* data [33], respectively.

ML-20M			Netflix			MSD			
models	nDCG @100	Recall @20	Recall @50	nDCG @100	Recall @20	Recall @50	nDCG @100	Recall @20	Recall @50
$\mathbf{\hat{B}}^{(\mathrm{dense})}$	0.423	0.392	0.522	0.397	0.364	0.448	0.391	0.334	0.430
MULT-VAE	0.426	0.395	0.537	0.386	0.351	0.444	0.316	0.266	0.364
MULT-DAE	0.419	0.387	0.524	0.380	0.344	0.438	0.313	0.266	0.363
CDAE	0.418	0.391	0.523	0.376	0.343	0.428	0.237	0.188	0.283
SLIM	0.401	0.370	0.495	0.379	0.347	0.428	-did not finish in [33]-		
WMF	0.386	0.360	0.498	0.351	0.316	0.404	0.257	0.211	0.312
training tim	es								
$\mathbf{\hat{B}}^{( ext{dense})}$	2 min 0 sec			1 min 30 sec			15 min 45 sec		
MULT-VAE	28 min 10 sec			1 hour 26 min			4 hours 30 min		
data-set properties	136,677 users 20,108 movies 10 million interactions			463,435 users 17,769 movies 57 million interactions			571,355 users 41,140 songs 34 million interactions		

## **EASE - the catch**



 EASE requires O(n\_items<sup>2</sup>) memory => <u>unfeasible</u> for large use-cases with millions of items





#### 2. How to scale them to huge datasets

## **ELSA**



- Assume that weight matrix B is a cosine similarity matrix
- B can be factorized  $B=AA^T$  (encoder = decoder)
- A is latent representation for items
- On the diagonal of B are ones
- Rows of A are L2-normalized
- Can be trained using numerical optimisation with backprop

## **ELSA**



- Weight matrix is factorized into low-rank structure
- Neither for optimizing the model nor predicting the recommendations, it will be necessary to store a dense matrix with dimensions superior to i×d



## **ELSA - results**



- Same performance as EASE
- Linear complexity with respect to the number of items



Comparison of EASE and ELSA: a) in terms of time complexity, b) memory complexity, c) performance on synthetic data





#### https://github.com/recombee/ELSA

## pip install elsarec





ELSA: memory requirements n\_items x latent\_dim (still large)

# Dense weight matrix of EASE can be accurately approximated by a product of 2 sparse matrices

- tool: factorized sparse approximate inversion
  - sophisticated numerical algorithm
  - can extract *global information* from the item-item network
  - factorization is the key to compression
- user chooses final model density

## SANSA

#### Architecture and training procedure

#### encoder decoder



Fig. 1. SANSA is a sparse nonsymmetric encoderdecoder model. To disallow recommending input items, we mask the prediction vector, or add an input-output residual connection. input user-item interaction matrix X, L2 regularization  $\lambda$ 1: compute sparse  $LDL^T \approx P(X^TX + \lambda I)P^T$  (for some permutation P) 2: compute sparse  $K \approx L^{-1}$ 3:  $W \leftarrow KP$ 4:  $Z_0 \leftarrow D^{-1}W$ 5:  $\vec{r} \leftarrow \text{diag}(W^TZ_0)$ 6:  $Z \leftarrow \text{scale the columns of } Z_0 \text{ by } -\vec{r}$ return  $W^T, Z$ 

Algorithm 1. The training procedure of SANSA is based on factorized sparse approximate inversion. The final scaling is applied to the decoder only.

## **SANSA - results**



#### Dense dataset



Sparse dataset

Amazon Books								
	results reprinted from [13]:							
	SANSA	MRF	MRF	EASER	SLIM	ITEM-	ULTRA-	
	(ICF)	(r=0)	(r = 0.5)			CF	GCN	
r@20	0.077	0.071	0.069	0.071	0.075	0.074	0.068	
n@20	0.064	0.058	0.055	0.057	0.060	0.061	0.056	
TRAINING RESOURCES								
vCPU	2	16	16	28	28	28	20*	
memory usage (GB):								
peak	9.18	96.45	96.58	—— not measured; costly ——				
avg.	3.87	49.12	49.75	—— not measured; costly ——				
time	49 s	172 s	167 s	222 m	316 m	57 m	45 m	



\*and a GPU (RTX 2080)

Figure 2: Accuracy of sansa (Cholmod) on MSD after various number of training scans s and finetune steps f.

Scaling hypothesis: larger dataset => greater compression At O(1M) items we get **100 000x compression** vs. full EASE





#### https://github.com/glami/sansa

### pip install sansa



#### 3. How to incorporate **item information**

## **ELSA** with item information



ELSA assumes weight matrix B is a cosine similarity matrix

Given images/descriptions/metadata, we can compute embeddings for items (SVD, deep learning...)

Idea:

- Instead of  $B=AA^T$  (encoder = decoder), let's represent B=ED
- Use the embedding matrix as either E or D
- Learn the other one using gradient descent

## ELSA with item information

Implications:

*Without retraining*, we can expand E or D by generating embeddings for items not seen during training (cold start), and

[expanded encoder] use new interactions with unseen items to recommend seen items Ε D X

 $\approx$ 



#### 4. How we teach autoencoders to generalize to new items

## beeFormer



Assume, that A is not initialized randomly, but computed from attributes (descriptions, images) with neural network (NN). Then, we can **propagate gradients from ELSA to NN that generated A** and update weights of that NN. However we face the following challenge:

In every training step, we need to encode all items in the catalogue (possibly millions of items) with NN. This is unfeasible for both time and GPU memory reason.

To make training feasible we employ three following tricks.



## beeFormer - trick #1 - leveraging sparsity



Interaction data are typically (very) sparse. Therefore, we can omit columns with only zeros (or limit them with negative sampling):



## beeFormer - trick #2 - gradient checkpointing 🔦

Storing the whole graph for backpropagation would be costly, therefore we split training into three steps:

- 1. Compute A without tracking gradients for NN
- 2. Perform ELSA step and track gradients for A (gradient checkpoint)
- 3. Compute A while tracking gradients to optimize NN.



## beeFormer - trick #3 - gradient accumulation

We don't need to compute gradients for all items at once during optimisation of NN.

We can do it in batches and accumulate gradients during the loop.

Finally we can update NN once with accumulated gradients.

Alg	orithm 1 beeFormer training step procedure in Python using PyTorch								
Inp	ut:								
	batch of interactions X								
Transformer model transformer									
descriptions of the selected items as sequences of tokens tokenized_texts									
	PyTorch optimizer optimizer optimizing the parameters of a Transformer model transformer								
Out	put:								
	Transformer model transformer with updated weights								
	predictions X_pred								
	computed loss loss								
1	with torch.no_grad():								
2	A_list = []								
3	<pre>for t in tokenized_texts:</pre>								

```
4 A_list.append(torch.nn.functional.normalize(transformer(t)))
```

```
5 A = torch.vstack(A_list)
```

```
o
7 A.requires_grad = True
```

```
8 X_pred = X @ A @ A.T - X
```

```
9
```

```
10 loss = torch.nn.MSELoss(
```

```
11 torch.nn.functional.normalize(X),
```

```
12 torch.nn.functional.normalize(X_pred),
```

```
13 )
14 loss.backward()
```

```
14 1055.backward()
15
```

```
16 checkpoint = A.grad
```

```
16 checkpoint = A.grad
17
```

```
17
18 optimizer.zero_grad()
```

```
19 for i, t in enumerate(tokenized_texts):
```

```
20 A_i = transformer(t)
```

```
21 A_i.backward(gradient=checkpoint[i])
```

```
22 optimizer.step()
```

## **beeFormer - results**

Dataset Method	Sentence Transformer	R@20	R@50	N@100					
GB10K CBF	Llama-goodbooks-mpnet Llama-goodbooks-nomic	0.2646 0.3070	0.3990 0.4248	0.3869 0.4161	Data Method	Transformer	R@20	R@50	N@100
GB10K Heater	all-mpnet-base-v2 nomic-embed-text-v1.5 bge-m3 Llama-goodbooks-mpnet Llama-goodbooks-nomic	$\begin{array}{r} 0.2078 \\ \underline{0.2154} \\ 0.2052 \\ 0.2570 \\ 0.2520 \end{array}$	$\begin{array}{c} 0.3221 \\ \underline{0.3317} \\ 0.3113 \\ 0.3800 \\ 0.3577 \end{array}$	$\begin{array}{r} \underline{0.3195}\\ 0.3193\\ 0.3099\\ 0.3749\\ 0.3644 \end{array}$	Images CBF	Img-amazfashion-clip-vit-b32 Img-amazfashion-clip-vit-114 Img-amazfashion-dinov2-small Img-amazfashion-dinov2-base Img-amazfashion-nomic	0.4216 <b>0.4490</b> 0.4169 0.4356 0.4403	0.4676 <b>0.4947</b> 0.4612 0.4739 0.4927	0.4071 <b>0.4262</b> 0.4024 0.4174 0.4148
ML20M CBF	Llama-movielens-mpnet Llama-movielens-nomic	0.4193 <b>0.4330</b>	<b>0.6147</b> 0.5404	<b>0.4066</b> 0.3874	Images	clip-vit-base-patch32 clip-vit-large-patch14	$0.3919 \\ 0.4076$	$0.4247 \\ 0.4407$	$0.3745 \\ 0.3903$
ML20M Heater	all-mpnet-base-v2 nomic-embed-text-v1.5 bge-m3	$\frac{0.3114}{0.3049}\\0.2847$	$\frac{0.4331}{0.4285}\\0.3932$	$\frac{0.3407}{0.3270}\\0.3161$	Heater	dinov2-small dinov2-base nomic-embed-vision-v1.5	$\begin{array}{c} 0.3914 \\ 0.3916 \\ \underline{0.4132} \end{array}$	$\begin{array}{c} 0.4264 \\ 0.4289 \\ \underline{0.4489} \end{array}$	0.3830 0.3828 <u>0.3994</u>
	Llama-movielens-mpnet Llama-movielens-nomic	$0.4320 \\ 0.4263$	$0.5979 \\ 0.5346$	$0.3937 \\ 0.3734$					

#### Dataset Sentence Transformer R@20 R@50 N@100 all-monet-base-v2 0.10170.1886 0.1739nomic-embed-text-v1.5 0.11460.20690.1896 bge-m3 0.11340.19530.1838 GB10K Llama-movielens-mpnet 0.17810.28560.2737Llama-amazbooks-mpnet 0.26480.39550.3792Llama-movielens-nomic 0.19640.2793 0.2832Llama-amazbooks-nomic 0.2908 0.4263 0.3998 all-mpnet-base-v2 0.1042 0.07880.1550nomic-embed-text-v1.5 0.1113 0.2143 0.1511 ML20M bge-m3 0.14090.21250.1578Llama-goodbooks-mpnet 0.16410.27710.2100 Llama-goodbooks-nomic 0.2710 0.40330.2695all-monet-base-v2 0.43860.4804 0.3497nomic-embed-text-v1.5 0.42340.47540.3480 Amazon bge-m3 0.44650.48820.3526fashion Llama-goodbooks-nomic 0.4800 0.52720.3995 Llama-movielens-nomic 0.46580.51110.3887

#### Cold start scenario - beeFormer trained models **generalize towards unobserved items**

Zero shot scenario beeFormer trained models transfer knowledge from one dataset to another



## **beeFormer - results**

Scenario Method	Model	R@20	R@50	N@100
	all-mpnet-base-v2	0.0218	0.0336	0.0193
	nomic-embed-text-v1.5	0.0387	0.0560	0.0320
	bge-m3	0.0398	0.0546	0.0313
zero-shot	Llama-goodbooks-mpnet	0.0650	0.0928	0.0513
CBF	Llama-goodbooks-nomic	0.0706	0.1005	0.0551
	Llama-movielens-mpnet	0.0266	0.0372	0.0221
	${\it Llama-movielens-nomic}$	0.0555	0.0734	0.0419
	KNN	0.0370	0.0562	0.0303
supervised	ALS MF	0.0344	0.0580	0.0313
CF	ELSA	0.0367	0.0628	0.0346
	SANSA	$\underline{0.0421}$	0.0678	0.0362
supervised	Llama-amazbooks-mpnet	0.0706	0.1045	0.0571
CBF	${\it Llama-amazbooks-nomic}$	0.0765	0.1150	0.0623

Dataset Scenario Model R@20 R@50 N@100 cold-start Llama-goodbooks-mpnet 0.26460.3990 0.3869 zero-shot Llama-movielens-mpnet 0.17810.28560.2737 0.3984cold-start Llama-goodlens-mpnet 0.26000.4109GB10K cold-start Llama-goodbooks-nomic 0.3070 0.42480.4161Llama-movielens-nomic 0.19640.28320.2793 zero-shot cold-start Llama-goodlens-nomic 0.28910.41130.4032zero-shot Llama-goodbooks-mpnet 0.16410.27710.2100cold-start Llama-movielens-mpnet 0.6147 0.4066 0.4193cold-start Llama-goodlens-mpnet 0.45190.6153 0.4106ML20M Llama-goodbooks-nomic 0.27100.40330.2695zero-shot cold-start Llama-movielens-nomic 0.4330 0.54040.3874cold-start Llama-goodlens-nomic 0.38990.49650.3557

Comparison with supervised CF models. Note that model trained on movies beats supervised CF models on books data. beeFormer can accumulate knowledge from multiple datasets (domains)







#### https://github.com/recombee/beeformer

#### https://huggingface.co/beeformer



### **Thank You For Your Attention**

## **Q & A**

