

SIEMENS

LMU

LUDWIG-  
MAXIMILIANS-  
UNIVERSITÄT  
MÜNCHEN

ESWC 2016 / Machine Learning Track

01.06.2016

# Embedding Mapping Approaches for Tensor Factorization and Knowledge Graph Modelling

*Yinchong Yang<sup>1,2</sup>, Cristóbal Esteban<sup>1,2</sup> and Volker Tresp<sup>1,2</sup>*

<sup>1</sup> Siemens AG, Corporate Technology, Munich, Germany

<sup>2</sup> Ludwig-Maximilians-Universität München, Munich, Germany

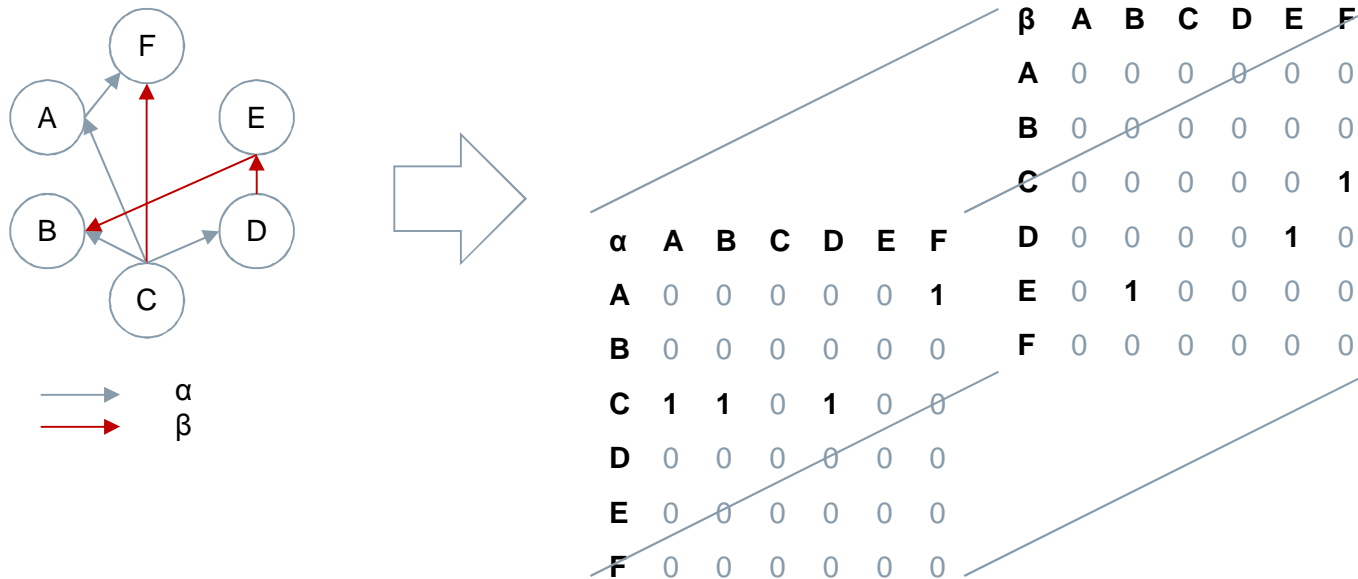
# Agenda

- ***Introduction to Latent Embedding Models***
- ***Motivation to Embedding Mapping***
- ***Embedding Mapping Approaches in Details***
- ***Experiments***
- ***Conclusions***

# Introduction to Latent Embedding Models

## The Starting Point: From Graph to Tensor

- Knowledge Graph represented as adjacency matrix/tensor  $\underline{Y}$ :
- $y_{i,j,k} = 1$  for the triple (entity  $i$ , relation  $k$ , entity  $j$ ) observed; and 0 otherwise;



$\underline{Y}$  of dimensions (6 x 6 x 2)

- 'Feature space' being usually highly dimensional and highly sparse;

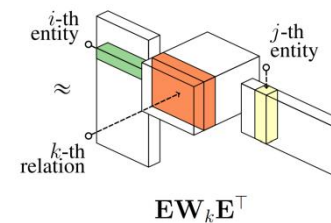
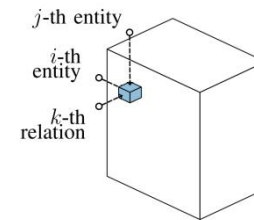
# Introduction to Latent Embedding Models

## Modelling of Adjacency Tensor

“Latent embedding models” aka. “factorization models” can model the KG-Tensors by

- Learning a latent embedding vector/matrix for each entity/relation;
- Learning an interaction rule of such vectors/matrices;
- And fulfil therefore the tasks i.a. of:
  - Prediction of unobserved links;
  - Distance-based entity resolution;
- Representative models: PARAFAC, Tucker, RESCAL, Multiway Neural Network, etc.;
- Successful applications in e.g.:
  - Recommender Systems [Koren et al. 2009] (matrices);
  - Knowledge Graph [Nickel et al. 2015, Dong et al. 2014] (tensors);

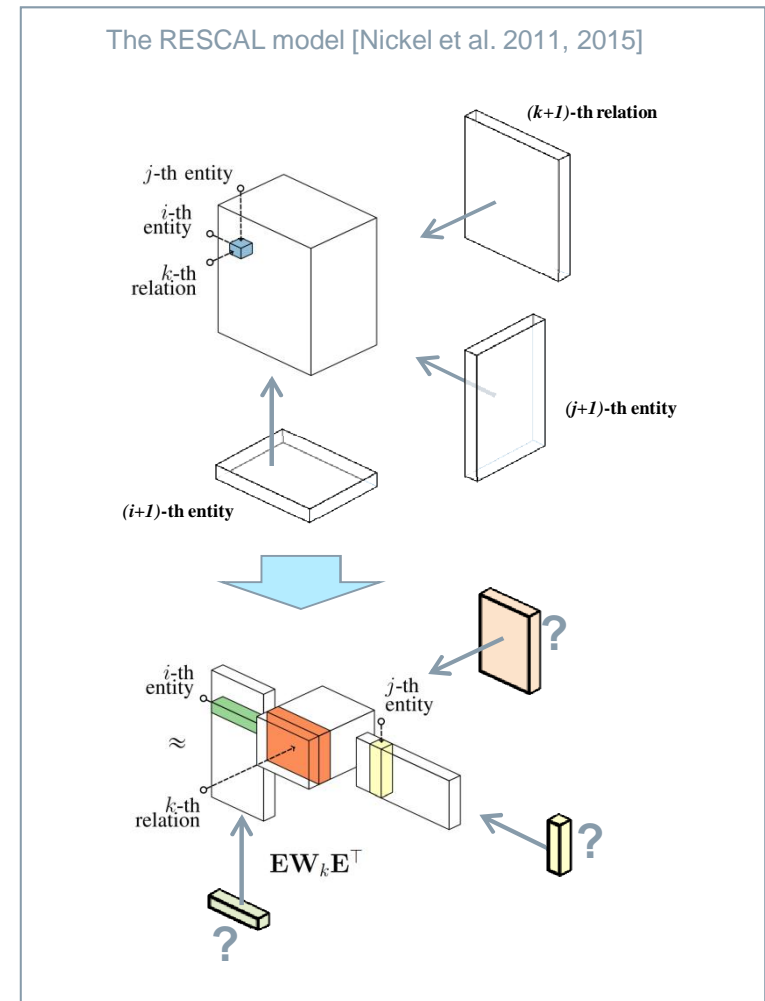
The RESCAL model [Nickel et al. 2011, 2015]



# Motivation to Embedding Mapping

## Embedding Mapping to handle new entities

- One limitation: deriving embeddings for new entities/relations ( $\rightarrow$ ):
  - Current known solution:
    - Retraining the entire model  $\rightarrow$  Not efficient;
    - KNN approximation  $\rightarrow$  Not effective;
- We introduce a Embedding Mapping ('Emma')
  - from the feature space (i.e. the set all known links between an entity and all other entities)
  - to the latent embedding space;



# Approach 0: Factorization Models with Closed-Form Mappings

## Canonical solutions

- There are a few factorization models where such mapping are for free as inverting. For instance:
  - SVD:  $X = UDV^T \rightarrow U = X(DV^T)^\dagger$  therefore generally assumed:  $u_{new}^T = x_{new}^T(DV^T)^\dagger$
  - Tucker:  $\mathcal{X} \approx \mathcal{G} \times_1 A \times_2 B \times_3 C$  i.e.  $X_{(1)} = AG_{(1)}(C \otimes B)^T$  analogously:  $a_{new}^T = x_{new}^T(G_{(1)}(C \otimes B)^T)^\dagger$
- What these models have in common:
  - Existence of matrix  $M_d$  independent of latent embeddings  $A^{(d)}$  for each dimension  $d$  so that:
    - $X_{(d)} = A^{(d)} (M_d)^\dagger$  i.e.  $A^{(d)} = X_{(d)} M_d$
- Other models do not enjoy such advantages:
  - RESCAL: shared embedding;
  - Multiway Neural Network (mwNN): non-linearity

# Approach 1: Post Mapping

*The most intuitive thought – “Emma-Post”*

- Step 1: Train the factorization model as usual;
- Step 2: Calculate a linear regression

$$\mathbf{a}_i^{(d)} = \mathbf{W}^{(d)T} \text{vec}(\mathbf{X}_i) + \mathbf{b}^{(d)},$$

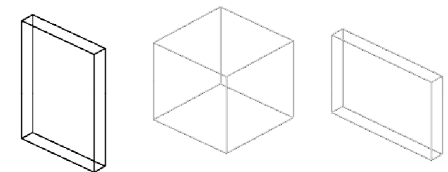
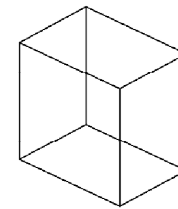
where  $\mathbf{a}_i^{(d)}$  is the embedding vector and  $\mathbf{X}_i$  is the  $i$ -th slice of the  $d$ -th dimension of the tensor  $\underline{\mathbf{X}}$ .

- So that the same regression relation generalizes to any new entity:

$$\mathbf{a}_{new}^{(d)} = \mathbf{W}^{(d)T} \text{vec}(\mathbf{X}_{new}) + \mathbf{b}^{(d)}$$

- Equivalence to inverting in cases of factorization models with closed-form mapping (details in paper);

Step 1:



# Approach 1: Post Mapping

*The most intuitive thought – “Emma-Post”*

- Step 1: Train the factorization model as usual;
- Step 2: Calculate a linear regression

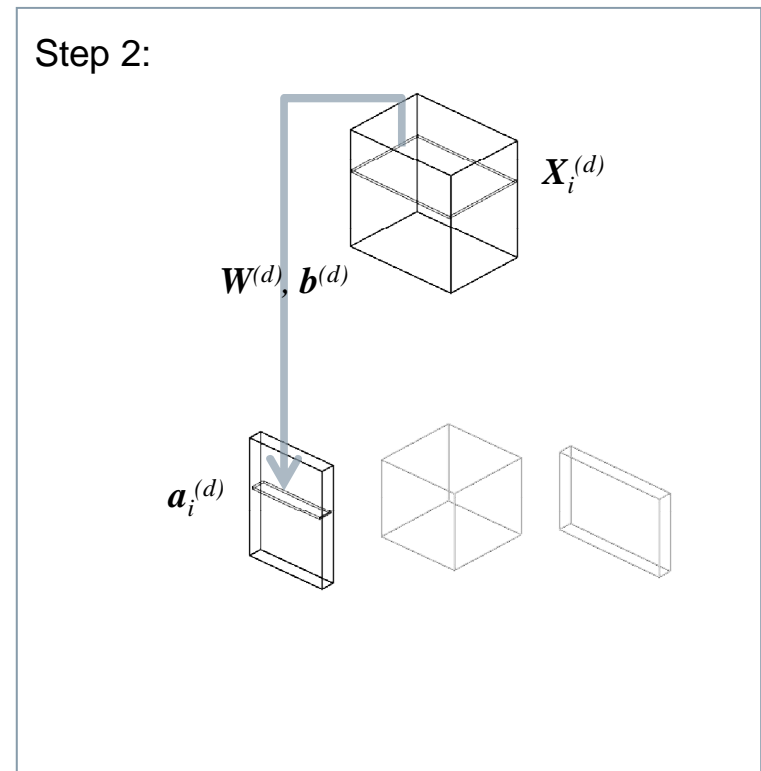
$$\mathbf{a}_i^{(d)} = \mathbf{W}^{(d)T} \text{vec}(\mathbf{X}_i) + \mathbf{b}^{(d)},$$

where  $\mathbf{a}_i^{(d)}$  is the embedding vector and  $\mathbf{X}_i$  is the  $i$ -th slice of the  $d$ -th dimension of the tensor  $\underline{\mathbf{X}}$ .

- So that the same regression relation generalizes to any new entity:

$$\mathbf{a}_{new}^{(d)} = \mathbf{W}^{(d)T} \text{vec}(\mathbf{X}_{new}) + \mathbf{b}^{(d)}$$

- Equivalence to inverting in cases of factorization models with closed-form mapping (details in paper);





# Approach 1: Post Mapping

*The most intuitive thought – “Emma-Post”*

- Step 1: Train the factorization model as usual;
- Step 2: Calculate a linear regression

$$\mathbf{a}_i^{(d)} = \mathbf{W}^{(d)T} \text{vec}(\mathbf{X}_i) + \mathbf{b}^{(d)},$$

where  $\mathbf{a}_i^{(d)}$  is the embedding vector and  $\mathbf{X}_i$  is the  $i$ -th slice of the  $d$ -th dimension of the tensor  $\underline{\mathbf{X}}$ .

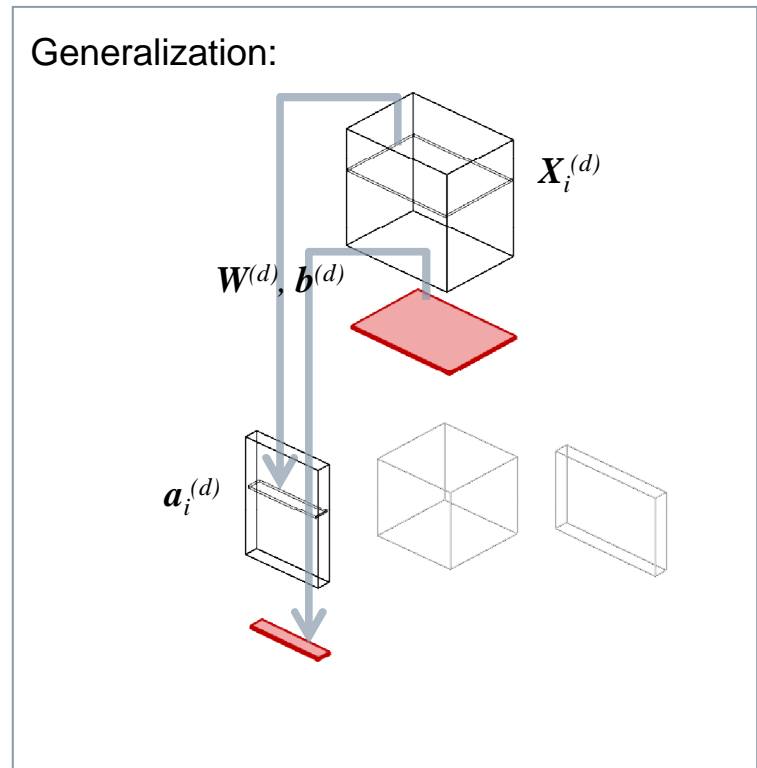
- So that the same regression relation generalizes to any new entity:

$$\mathbf{a}_{new}^{(d)} = \mathbf{W}^{(d)T} \text{vec}(\mathbf{X}_{new}) + \mathbf{b}^{(d)}$$

- Equivalence to inverting in cases of factorization models with closed-form mapping (details in paper);

Reducing mapping error function after factorization error;

- + Simple, applicable for arbitrary factorization model;
- Regression error reduced only to a certain extent



## Approach 2: Hatting Algorithm

*Integrating the linear regression into the factorization – “Emma-Hatting”*

- Calculate the linear regression after each iteration (ALS, gradient approaches) of the factorization learning (instead of only once after factorization completes);
- → Replacing the latent embedding vector with its Least-Square estimates (therefore ‘Hatting’) in each iteration:

```

for each iteration of the factorization:
  for each dimension  $d$  of target tensor  $\underline{X}$ :
    calculate the latent embedding  $\mathbf{A}^{(d)}$  as prescribed by the factorization model
    replace  $\mathbf{A}^{(d)}$  with its LS estimate w.r.t.  $\mathbf{X}_{(d)}$  as regressor:
       $\mathbf{A}^{(d)} \leftarrow \hat{\mathbf{A}}_{LS}^{(d)} = \mathbf{X}_{(d)} (\mathbf{X}_{(d)}^T \mathbf{X}_{(d)} - \lambda \mathbf{I})^{-1} \mathbf{X}_{(d)}^T \mathbf{A}^{(d)}$ 
  return  $\mathbf{A}^{(d)}$  as latent embeddings;
   $\mathbf{M}_d = (\mathbf{X}_{(d)}^T \mathbf{X}_{(d)})^{-1} \mathbf{X}_{(d)}^T \mathbf{A}^{(d)}$  as mapping matrices
  for each  $d$ 

```

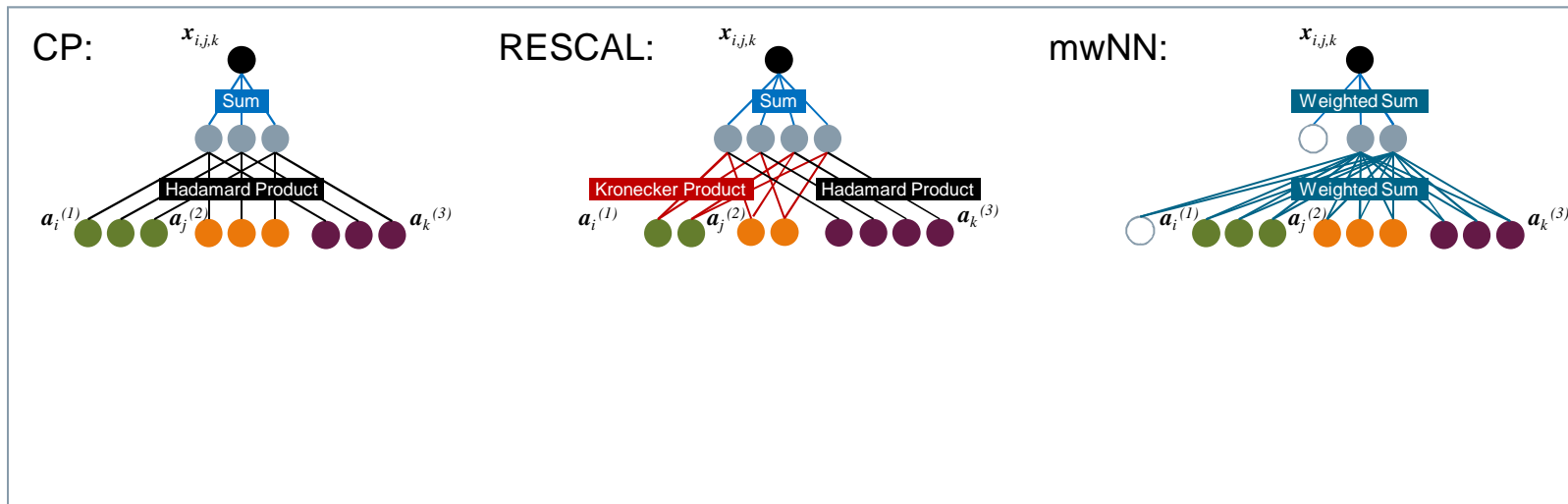
Reducing the two errors of factorization and mapping alternatively;

- + The regression error almost eliminated;
- + Though applied in each iteration, the update consists of only one matrix multiplication;
- But how does the Hatting update affects the gradient approaches?

# Approach 3: Back-Propagation

*Inspired by an Neural Network aspect – “Emma-BP”*

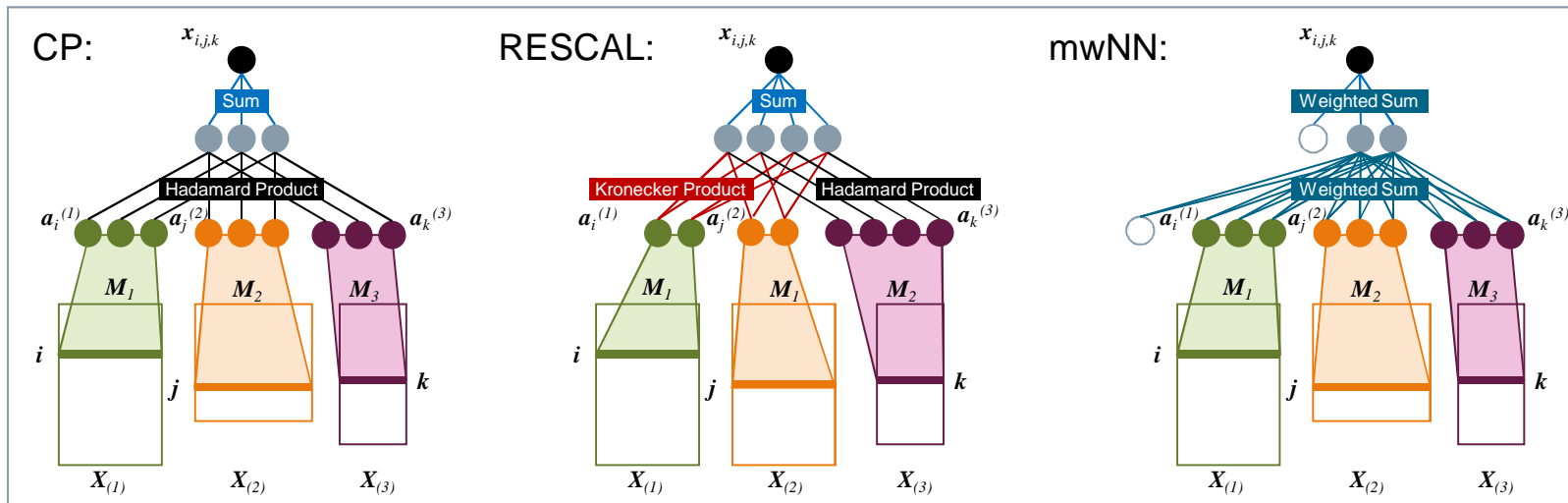
- Most factorization models rewritable as Neural Networks:
- 1) input layer = latent embeddings; 2) output = tensor entry; 3) architecture = the interaction rule:
  - $\mathbf{x}_{i,j,k} = f_{NN}(\mathbf{a}_i^{(1)}, \mathbf{a}_j^{(2)}, \mathbf{a}_k^{(3)})$



# Approach 3: Back-Propagation

*Inspired by an Neural Network aspect – “Emma-BP”*

- Most factorization models rewritable as Neural Networks:
- 1) input layer = latent embeddings; 2) output = tensor entry; 3) architecture = the interaction rule:
  - $\mathbf{x}_{i,j,k} = f_{NN}(\mathbf{a}_i^{(1)}, \mathbf{a}_j^{(2)}, \mathbf{a}_k^{(3)})$
- Emma as one more linear layer  $\rightarrow$  train factorization + Emma models as a compact one with BP
  - $\mathbf{x}_{i,j,k} = h_{NN}(\mathbf{x}_{(1)}_i, \mathbf{x}_{(2)}_j, \mathbf{x}_{(3)}_k)$

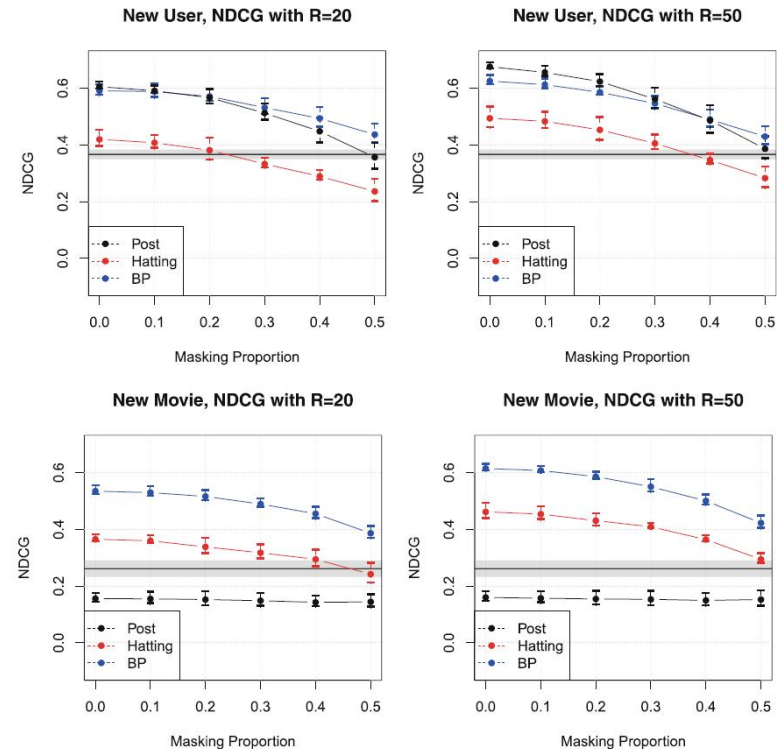


Reducing one compact cost function;  
 + Guarantees a perfect mapping i.e. without mapping error;  
 - Numerically more expensive;

# Experiment 1: MovieLens Data

*Emma applied for matrix factorization and Collaborate Filtering*

- **Data:** a binary user-item matrix (943 x 1682);
- **Task:** make movie recommendations to new users and user recommendations to new movies;
- **Settings:**
  - Baseline: Most-frequent Prediction;
  - Sampling 20% users/movies as test set;
  - Masking a proportion of [0%, 10%, 20%, 30%, 40%, 50%] of all entries, simulating incomplete information;
  - Mapping the test set into the latent space and calculate recommendations (recovering the masked ones)
  - Evaluation in term of NDCG@k
- **Models:** Matrix Factorization according to [Koren et al. 2009] x all 3 Emma approaches
- **Results:**
  - New users: Emma-Post, Emma-BP > Emma-Hatting
  - New movies: Emma-BP > Emma-Hatting > Emma-Post



*Use case: Real-time recommendation for new user/item  
Further reference: Rendle's MyMediaLite\* package  
provides NDCG's between 0.56 and 0.62, when  
generating recommendations for known users.*

\*

[http://www.mymedialite.net/examples/item\\_recommendationDatasets.html](http://www.mymedialite.net/examples/item_recommendationDatasets.html) Siemens AG 2016. All rights reserved

## Experiment 2: FreeBase Data

*Emma applied for tensor factorization and KG link prediction*

- **Data:** A sampled fraction of KG, represented as a binary tensor (39 x 115 x 115);
- **Task:** Predicting links for new entities
- **Settings:** 'Baseline': Retraining;
  - Sampling 20% of all entities and masking 20% links;
  - Mapping the entities into the latent space and to predict the masked links;
  - Evaluation in term of AUROC, AUPRC
- **Models:** RESCAL, mwNN in combination with all 3 Emma approaches
- **Results:**
  - mwNN outperform RESCAL, confirming [Krompaß et al. 2015];
  - Emma-Hatting and Emma-BP superior to Emma-Post;
  - The combination of mwNN and Emma-BP show performances close even to those from retraining;

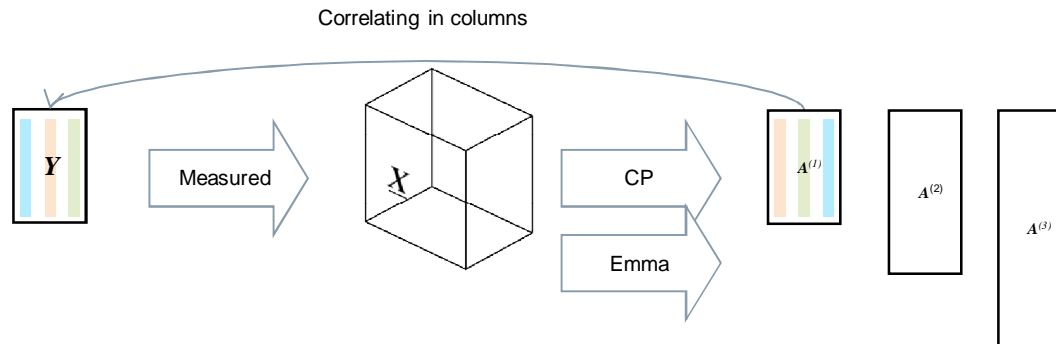
Fac.	Mapping	R = 10		R = 30	
		AUROC	AUPRC	AUROC	AUPRC
RESCAL	Retraining	0.901 ± 0.039	0.820 ± 0.059	0.788 ± 0.054	0.616 ± 0.100
	Emma-Post	0.759 ± 0.096	0.600 ± 0.145	0.693 ± <b>0.065</b>	0.432 ± <b>0.106</b>
	Emma-Hatting	<b>0.778</b> ± 0.112	<b>0.605</b> ± 0.152	0.700 ± 0.091	0.481 ± 0.120
	Emma-BP	0.700 ± <b>0.060</b>	0.485 ± <b>0.092</b>	<b>0.740</b> ± 0.090	<b>0.509</b> ± 0.134
mwNN	Retraining	0.964 ± 0.008	0.886 ± 0.060	0.970 ± 0.010	0.923 ± 0.017
	Emma-Post	0.844 ± <b>0.009</b>	0.390 ± 0.101	0.826 ± <b>0.035</b>	0.382 ± <b>0.063</b>
	Emma-Hatting	0.847 ± 0.042	0.423 ± 0.118	0.843 ± 0.038	0.394 ± 0.081
	Emma-BP	<b>0.949</b> ± 0.022	<b>0.805</b> ± <b>0.080</b>	<b>0.931</b> ± 0.036	<b>0.735</b> ± 0.101

# Experiment 3: Amino Acid Data

*Verifying Emma's plausibility and interpretability of the latent space*

- **Data:**

- 3 types of amino acid mixed according to 5 recipes: a real matrix of (5 x 3)  $Y$ ;
- The 5 samples measured by fluorescence with excitation 250-300nm, emission 250-450nm: a real tensor of (5 x 61 x 201)  $X$ ;
- With a rank 3 CP factorization, the latent embeddings representing recipes ( $A^{(1)}$ ) are expected to correlate column-wise perfectly with the recipe matrix itself (whilst the column order may vary).



- **Task:**

- What if there is a new measurement matrix (61 x 201) of a mixture according to a new recipe? Will the correspondingly mapped embedding vector still correlate as before?

- **Settings:** Leave-1-out and Leave-2-out, evaluation in term of Pearson Correlation

- **Model:** CP in combination with Emma-Post and Emma-Hatting

- **Results:** Leave-1-out: both 0.999; Leave-2-out: both 0.991;

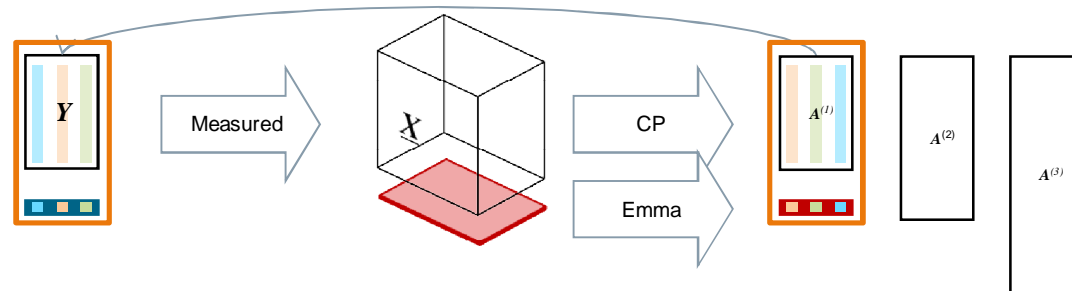
# Experiment 3: Amino Acid Data

*Verifying Emma's plausibility and interpretability of the latent space*

- **Data:**

- 3 types of amino acid mixed according to 5 recipes: a real matrix of (5 x 3)  $Y$ ;
- The 5 samples measured by fluorescence with excitation 250-300nm, emission 250-450nm: a real tensor of (5 x 61 x 201)  $X$ ;
- With a rank 3 CP factorization, the latent embeddings representing recipes ( $A^{(1)}$ ) are expected to correlate column-wise perfectly with the recipe matrix itself (whilst the column order may vary).

Correlating in columns still ?



- **Task:**

- What if there is a new measurement matrix (61 x 201) of a mixture according to a new recipe? Will the correspondingly mapped embedding vector still correlate as before?

- **Settings:** Leave-1-out and Leave-2-out, evaluation in term of Pearson Correlation

- **Model:** CP in combination with Emma-Post and Emma-Hatting

- **Results:** Leave-1-out: both 0.999; Leave-2-out: both 0.991;



# Conclusions and Outlooks

- Embedding Mapping Approaches derive latent embeddings for new entities in real time;
- Two classes of mapping:
  - a) Explicit mapping that are learned
    - after the factorization (Emma-Post);
    - during the factorization (Emma-Hatting);
  - b) Implicit mapping that are learned
    - as one linear activated layer of specific NN model (Emma-BP);
    - as closed-form mapping (inverting) in case of certain models;
- Experiments
  - a) To predict new links in recommender system and knowledge graph;
  - b) To verify interpretability of mapped embeddings;
- Future works
  - a) Improve the scalability by attacking the bottleneck due to unfolding operation;
  - b) Combining the work of [Gantner et al., 2010]: two-fold mapping: content attributes → feature space → latent space

## (Selected) References

- Nickel, Maximilian, Volker Tresp, and Hans-Peter Kriegel. "A three-way model for collective learning on multi-relational data." *Proceedings of the 28th international conference on machine learning (ICML-11)*. 2011.
- Nickel, Maximilian, et al. "A review of relational machine learning for knowledge graphs: From multi-relational link prediction to automated knowledge graph construction." *arXiv preprint arXiv:1503.00759* (2015).
- Krompaß, Denis, Stephan Baier, and Volker Tresp. "Type-constrained representation learning in knowledge graphs." *The Semantic Web-ISWC2015*. Springer International Publishing, 2015. 640-655.
- Koren, Yehuda, Robert Bell, and Chris Volinsky. "Matrix factorization techniques for recommender systems." *Computer* 8 (2009): 30-37.
- Gantner, Zeno, et al. "Learning attribute-to-feature mappings for cold-start recommendations." *Data Mining (ICDM), 2010 IEEE 10th International Conference on*. IEEE, 2010.
- Dong, Xin, et al. "Knowledge vault: A web-scale approach to probabilistic knowledge fusion." *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2014.

*Vielen Dank für Ihre Aufmerksamkeit*

*Thank you for your attention*

*πολλές ευχαριστίες*

**Yinchong Yang**

Doktorand

Ludwig-Maximilians-Universität München

Siemens AG, Corporate Technology

Kontakt:

[yinchong.yang@siemens.com](mailto:yinchong.yang@siemens.com)