

Short Bio

2009-2015: Bachelor and Master of statistics, minor in computer science. LMU Munich.

2014-2018: Dr. rer. nat. in computer science. LMU Munich, supervised by Prof. Dr. Volker Tresp.

2018-2022: Research Scientist at Siemens AG, Munich.

2022- : Senior Key Expert for robust AI at Siemens AG, Munich.

2019- : Co-lecturer “Deep Learning and AI” at LMU Munich.

Tensorization and uncertainty quantification in Deep Learning

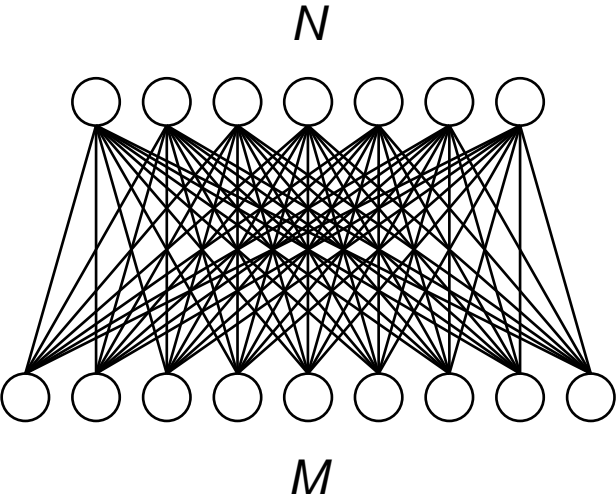
Dr. Yinchong Yang, Siemens AG

Invited Talk at Cyprus Institute, 2023-05-23

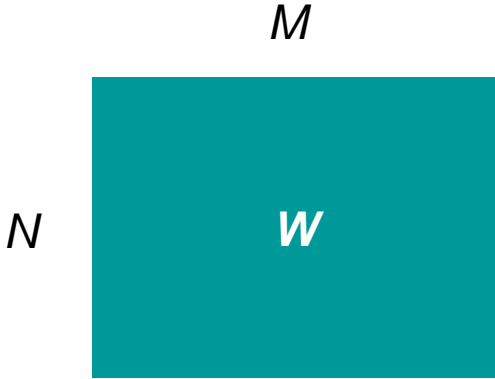
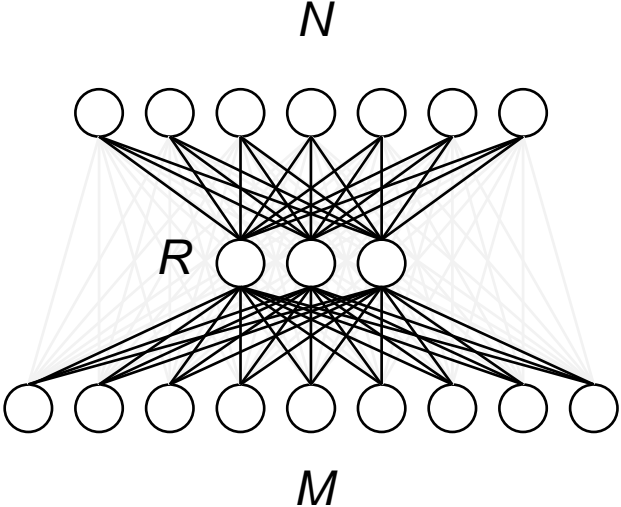
Tensorization of neural networks in Deep Learning

“Representing large weight matrices in NN with product of small tensors...”

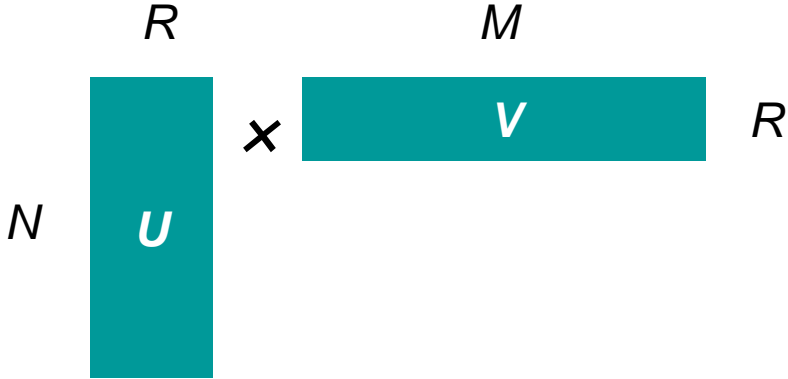
A motivating example



\approx



\approx



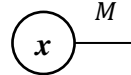
Preliminary: Tensor Graph and Einstein Summation Notation

In the syntax of tensor graph notation:

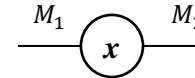
- Node = tensor
- Edge = “mode” / dimensionality



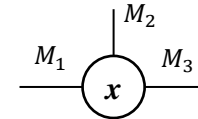
$$x \in \mathbb{R}$$



$$x \in \mathbb{R}^M$$



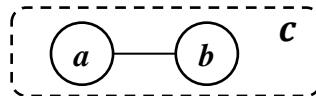
$$X \in \mathbb{R}^{M_1 \times M_2}$$



$$X \in \mathbb{R}^{M_1 \times M_2 \times M_2}$$

- Connection = tensor contraction
 - Combining multiple tensors to form a new one.
 - A graph representation of Einstein’s summation.

Tensor graph notation



Matrix notation

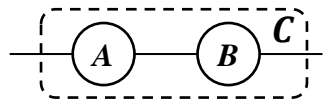
$$c = \langle a, b \rangle$$

Einstein summation notation

$$c = a(n)b(n)$$

pytorch

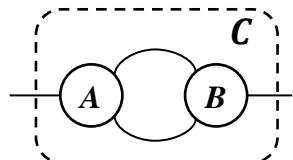
```
a = torch.randn(5)
b = torch.randn(5)
torch.einsum('n,n->', a, b)
# torch.inner(a, b)
```



$$C = AB$$

$$C(m, k) = A(m, n)B(n, k)$$

```
A = torch.randn(3, 5)
B = torch.randn(5, 4)
torch.einsum('mn,nk->mk', A, B)
# torch.matmul(A, B)
```



?

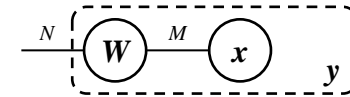
$$C(n, r) = A(m, n, k)B(m, k, r)$$

```
A = torch.randn(5, 10, 7)
B = torch.randn(5, 7, 2)
torch.einsum('mnk,mkr->nr', A, B)
```

Tensorized neural networks

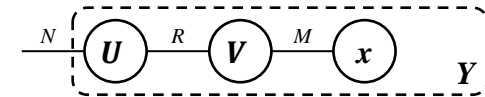
- A fully connected layer (omitting the bias) $\mathbf{y} = \mathbf{W}\mathbf{x}$ can be formulated as contraction:

- $\mathbf{y}(n) = \mathbf{W}(n, m)\mathbf{x}(m)$, with $\mathbf{W} \in \mathbb{R}^{N \times M}$, $\mathbf{x} \in \mathbb{R}^M$.
- Size: $N \times M$.



- A bottleneck FC layer as:

- $\mathbf{y}(n) = \mathbf{V}(n, r)\mathbf{U}(r, m)\mathbf{x}(m)$ with $\mathbf{V} \in \mathbb{R}^{N \times R}$, $\mathbf{U} \in \mathbb{R}^{R \times M}$, $\mathbf{x} \in \mathbb{R}^M$.
- Size: $R \times (M + N)$.

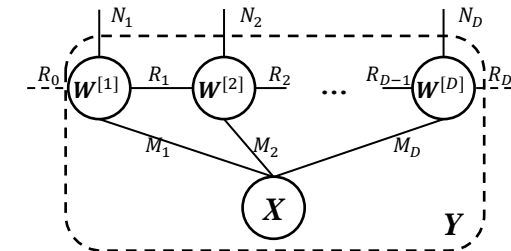
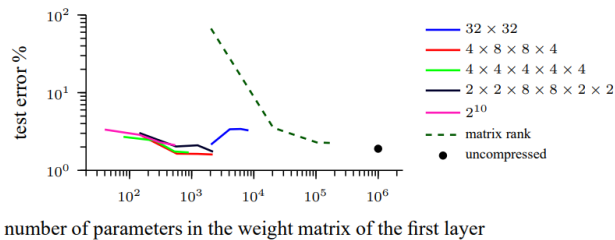


- A tensor-train layer [Novikov et al. 2015] closely related to Matrix Product State.

- $\mathbf{y}(n_1, n_2, \dots, n_D) = \mathbf{W}^{[D]}(n_D, r_{D-1}, r_D, m_D) \dots \mathbf{W}^{[2]}(n_2, r_1, r_2, m_2) \mathbf{W}^{[1]}(n_1, r_0, r_1, m_1) \mathbf{x}(m_1, m_2, \dots, m_D)$,
- Size: $\sum_{d=1}^D M_d N_d R_{d-1} R_d$ instead of $M \times N = \prod_{d=1}^D M_d N_d$

$$M = M_1 \times M_2 \times \dots \times M_D$$

$$N = N_1 \times N_2 \times \dots \times N_D$$



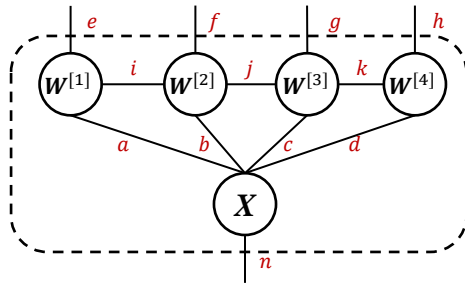
Novikov, Alexander, et al. "Tensorizing neural networks." *NeurIPS* 2015.

Tensorized neural networks

- Implementation in pytorch:
 - Step 1: decide the decomposition of input and output modes, as well as the rank.
 - Example: $4096 \rightarrow 16$ can be decomposed into $(8 \times 8 \times 8 \times 8) \rightarrow (2 \times 2 \times 2 \times 2)$ with rank 2

- Step 2: draw a tensor train layer and name the edges.

- Example:



```
opt = torch.optim.Adam([W1, W2, W3, W4], lr=5e-2)
loss = torch.nn.MSELoss()

for i in range(10000):
    y_hat = torch.einsum('aei, bfij, cgjk, dhk, nabcd -> nefgh', W1, W2, W3, W4, X)
    error = loss(y_hat.reshape(n_samples, -1), y.reshape(n_samples, -1))
    error.backward()
    opt.step()
    opt.zero_grad()
```

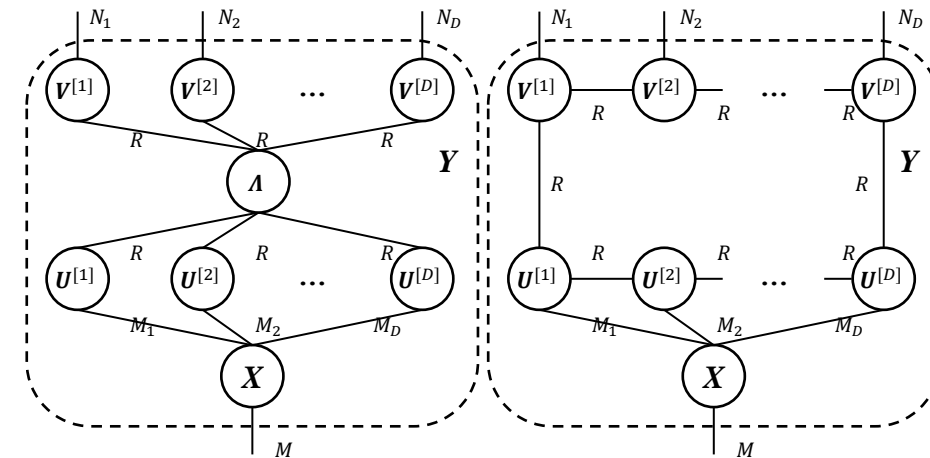
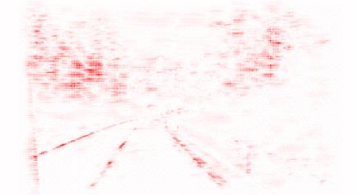
- Step 3: implement an torch.einsum operation with given names in step 2.
- Step 4: implement the backward pass.

Tensorized neural networks

	Accuracy	# Parameters	Runtime
TT-MLP	0.427 ± 0.045	7,680	902s
GRU	0.488 ± 0.033	44,236,800	7,056s
LSTM	0.492 ± 0.026	58,982,400	8,892s
TT-GRU	0.813 ± 0.011	3,232	1,872s
TT-LSTM	0.796 ± 0.035	3,360	2,160s

Original: (Liu et al., 2009)	0.712
(Liu et al., 2013)	0.761
(Hasan & Roy-Chowdhury, 2014)	0.690
(Sharma et al., 2015)	0.850
Our best model (TT-GRU)	0.813

- Line of research: new architectures based on tensorized neural network:
 - Tensorized CNN [Garipov et al. 2016],
 - Tensorized RNN [Yang et al. 2017, Tjandra et al. 2017],
- Line of research: new applications:
 - Time series data [Yu et al. 2017]
 - Video data – e.g. autonomous driving (contribution by Max Pittner) with LRP.
 - Sparse tabular data [Yang et al. 2017]
- Line of research: other decompositions
 - CP and Tucker [Pan et al. 2023] (left)
 - Tensor-Ring [Wang et al. 2018] (right)



Garipov, Timur, et al. "Ultimate tensorization: compressing convolutional and fc layers alike." *NeurIPS (2016) Workshop*.

Tjandra, Andros, Sakriani Sakti, and Satoshi Nakamura. "Compressing recurrent neural network with tensor train." *2017 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2017.

Yu, Rose, et al. "Long-term forecasting using tensor-train rnns." *Arxiv* (2017).

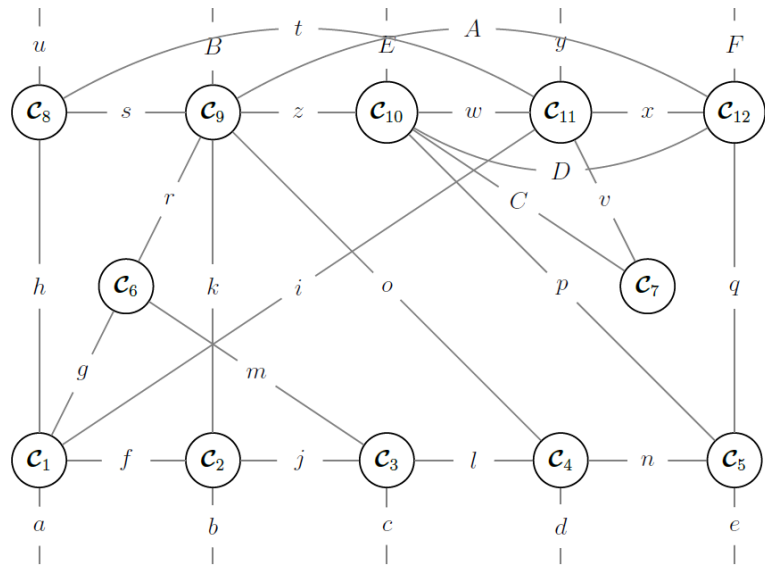
Yang, et al. "Modeling progression free survival in breast cancer with tensorized recurrent neural networks and accelerated failure time models." *MLHC, PMLR*, 2017.

Pan et al. <https://tednet.readthedocs.io/en/latest/> visited on 2023-05-18

Wang, Wenqi, et al. "Wide compression: Tensor ring nets." *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018.

Tensorized neural networks

- Flexible topology [Li and Sun 2020, Master thesis Arber Qoku]
 - Representing the core tensor topology with adjacency matrix (vectorized lower-/upper triangle).
 - Optimize the topology of core tensors as a hyper-parameter via evolution strategy and Bayesian optimization:



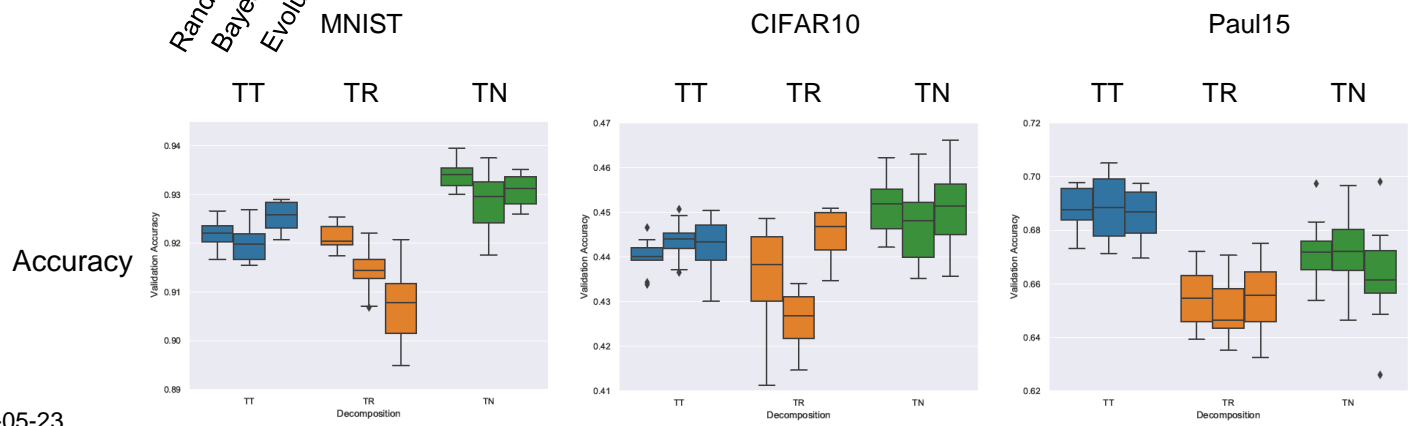
Li, Chao, and Zhun Sun. Evolutionary topology search for tensor network decomposition. *ICML*, 2020.

$$A_{TT} = \begin{pmatrix} 0 & R & 0 & 0 \\ R & 0 & R & 0 \\ 0 & R & 0 & R \\ 0 & 0 & R & 0 \end{pmatrix} \equiv \begin{matrix} \textcircled{C_1} & R & \textcircled{C_2} & R & \textcircled{C_3} & R & \textcircled{C_4} \\ i_1 & & i_2 & & i_3 & & i_4 \end{matrix}$$

$$A_{TR} = \begin{pmatrix} 0 & R & 0 & R & 0 \\ R & 0 & R & 0 & 0 \\ 0 & R & 0 & 0 & R \\ R & 0 & 0 & 0 & R \\ 0 & 0 & R & R & 0 \end{pmatrix} \equiv \begin{matrix} \textcircled{C_1} & R & \textcircled{C_2} & R & \textcircled{C_3} \\ R & & & & R \\ \textcircled{C_4} & & & R & \textcircled{C_5} \\ i_4 & & & i_5 & \\ i_4 & & i_5 & & i_6 \end{matrix}$$

$$A_{TN} = \begin{pmatrix} 0 & R & 0 & R & 0 & R & 0 & 0 \\ R & 0 & R & 0 & R & 0 & 0 & 0 \\ 0 & R & 0 & 0 & R & 0 & 0 & 0 \\ R & 0 & 0 & 0 & R & 0 & R & 0 \\ 0 & R & R & R & 0 & 0 & 0 & R \\ R & 0 & 0 & 0 & 0 & 0 & R & 0 \\ 0 & 0 & 0 & R & 0 & R & 0 & R \\ 0 & 0 & 0 & 0 & R & 0 & R & 0 \end{pmatrix} \equiv \begin{matrix} \textcircled{C_6} & R & \textcircled{C_7} & R & \textcircled{C_8} \\ R & & & & R \\ R & \textcircled{C_4} & & R & \textcircled{C_5} \\ R & & & R & R \\ \textcircled{C_1} & R & \textcircled{C_2} & R & \textcircled{C_3} \\ i_1 & & i_2 & & i_3 \end{matrix}$$

Random Search
Bayesian optimization
Evolution strategy



Uncertainty quantification with GP in Deep Learning

“Combining the predictive distribution of GP with representation learning by NN.”

Gaussian Processes Recap

- Standard Gaussian Process regression:

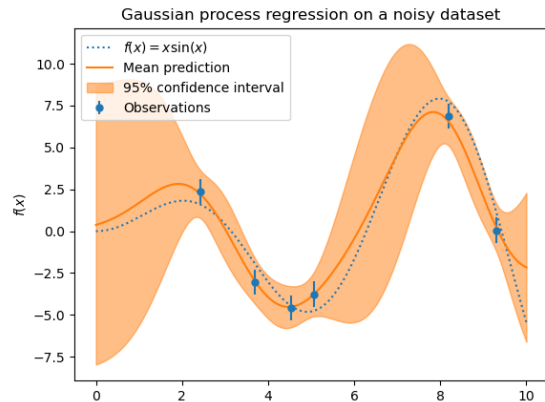
- $$p(y_* | \mathbf{x}_*, D) = \int p(y_* | \mathbf{x}_*, \boldsymbol{\theta}) p(\boldsymbol{\theta} | D) d\boldsymbol{\theta}$$

$$= N\left(\frac{1}{\sigma^2} \mathbf{k}_*^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}, k_{**} - \mathbf{k}_*^T (\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{k}_* + \sigma^2\right)$$

with

$$\mathbf{k}_* = (k(\mathbf{x}_*, \mathbf{x}_i))_{i=1}^n \in \mathbb{R}^n \text{ and } \mathbf{K} = (k(\mathbf{x}_i, \mathbf{x}_j))_{i,j=1}^{n,n} \in \mathbb{R}^{n \times n}$$

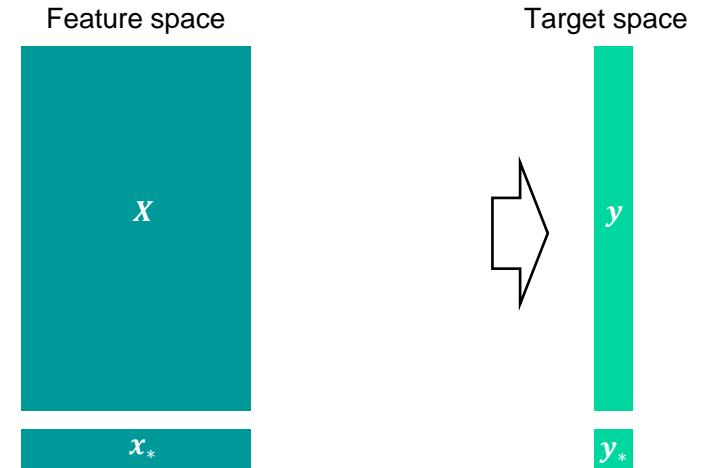
being the kernel function that measures the similarity between data points.



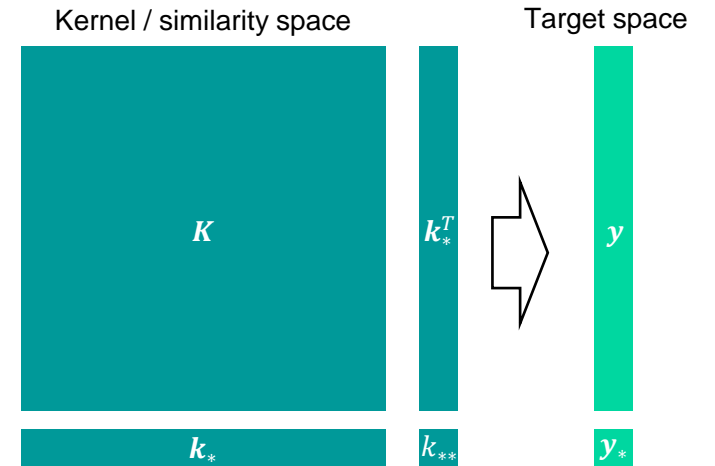
https://scikit-learn.org/stable/auto_examples/gaussian_process/plot_gpr_noisy_targets.html

- The computational complexity is $\mathcal{O}(n^3)$ due to the matrix inverse.

Parametric models:



Non-Parametric models:



Scalable Gaussian Processes

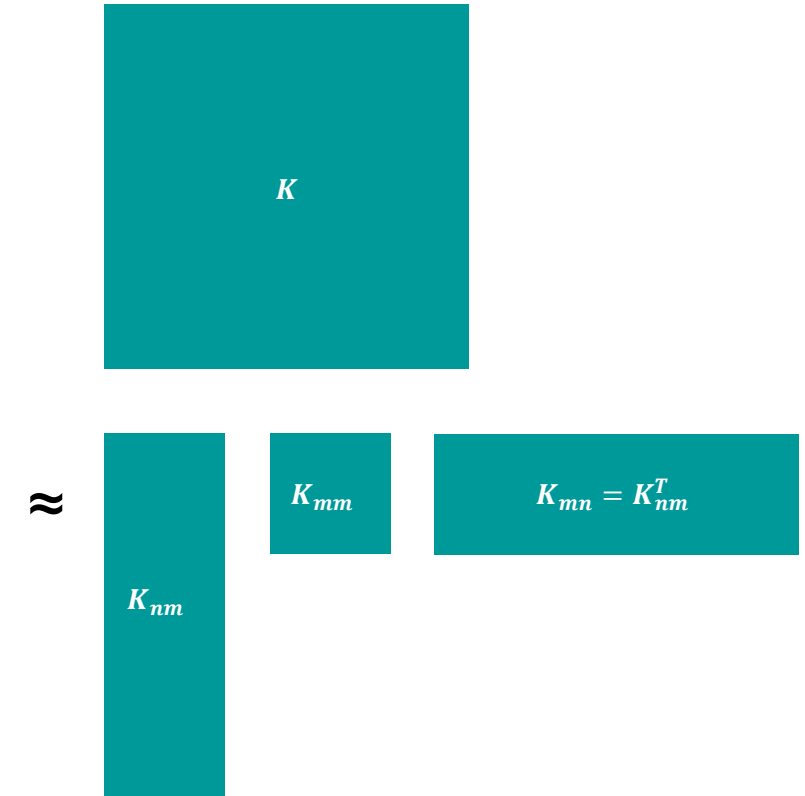
- Various approaches have been proposed to approximate the full covariance matrix. But the idea remains the same: representing the whole training set with a few “*inducing points*”:

$$\mathbf{K} \approx \mathbf{Q} = \mathbf{K}_{nm} \mathbf{K}_{mm}^{-1} \mathbf{K}_{mn}$$

so that it can be inverted using the Woodbury formula [Williams & Seeger 2000]. That is,

$$k(\mathbf{x}_i, \mathbf{x}_j) \approx \mathbf{k}(\mathbf{x}_i, \mathbf{Z})^T \mathbf{K}_{mm}^{-1} \mathbf{k}(\mathbf{x}_j, \mathbf{Z})$$

- The set of inducing points \mathbf{Z} could also be learned as parameter:
 - FITC [Snelson et al. 2000]
 - VFE [Titsias 2009]
 - SVGP [Hensman et al. 2013]
 - PPGP [Jankowiak et al. 2020]



Williams and Seeger. Using the Nyström method to speed up kernel machines. *NeurIPS* 2000.

Snelson et al. Sparse Gaussian processes using pseudo-inputs. *NeurIPS*. 2006.

Titsias. Variational learning of inducing variables in sparse Gaussian processes. *AISTATS* 2009.

Hensman et al. Gaussian Processes for Big Data. *UAI*. 2013.

Jankowiak et al. Parametric gaussian process regressors. *ICML*, 2020.

Scalable Gaussian Processes for collaborative filtering

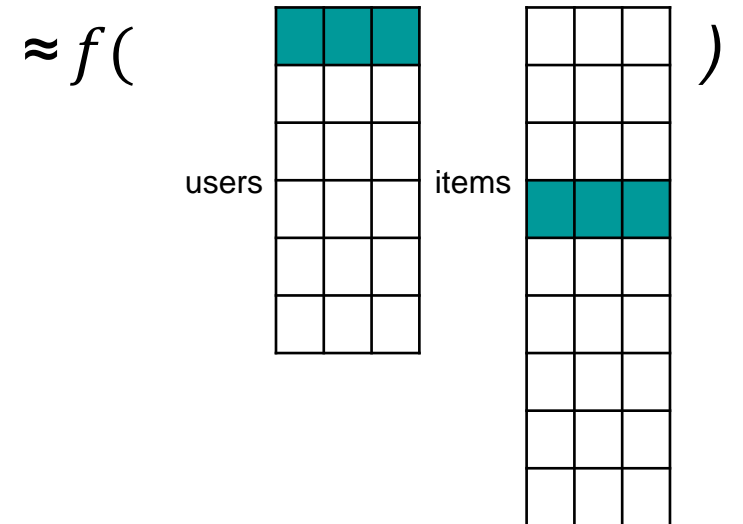
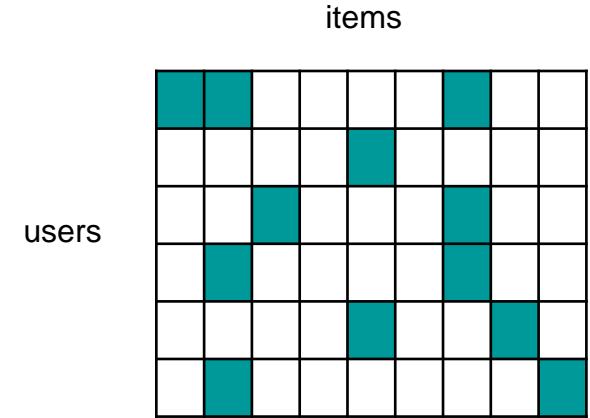
- Collaborative filtering: modeling user-item interaction, applied in recommender systems.
- Classical approaches include matrix decomposition such as SVD and NMF, more advanced methods such as multiway NN [Nickel et al. 2015].

$$y_{i,j} = f(\mathbf{a}_i, \mathbf{b}_j)$$

- with

$$\mathbf{Y} \in \mathbb{R}^{I \times J}, \mathbf{A} \in \mathbb{R}^{I \times r}, \mathbf{B} \in \mathbb{R}^{I \times r}$$

- The expressiveness of the model depends on the choice of function f .
- GP-LVM is one of the very few methods are capable of expressing predictive uncertainty, but doesn't scale well to large and sparse user-item matrices.

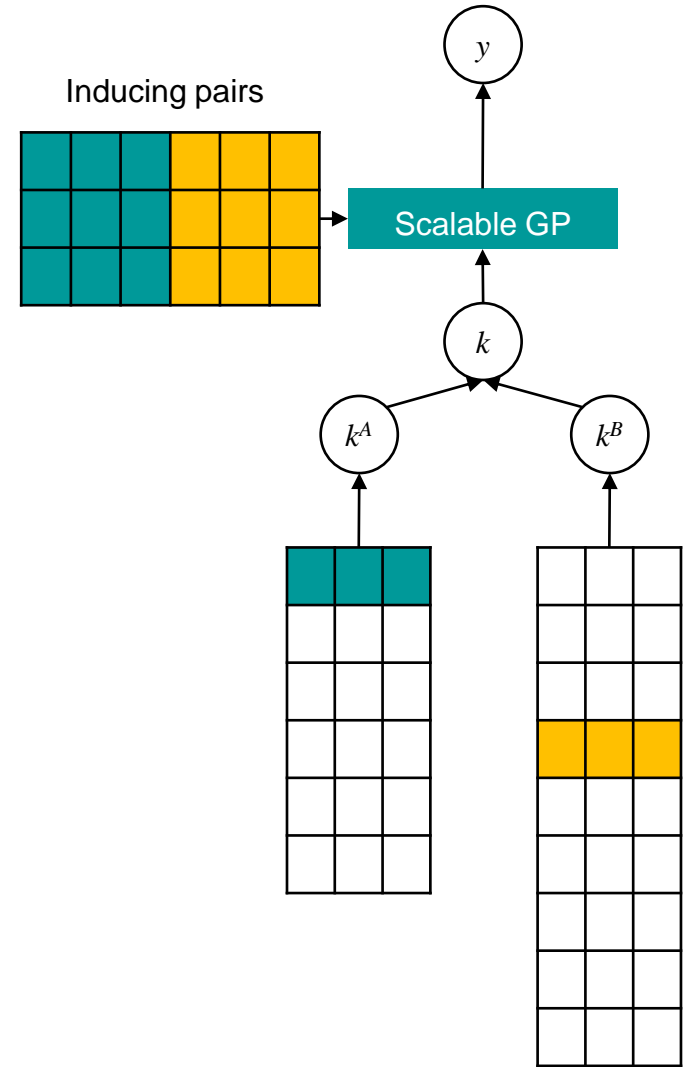
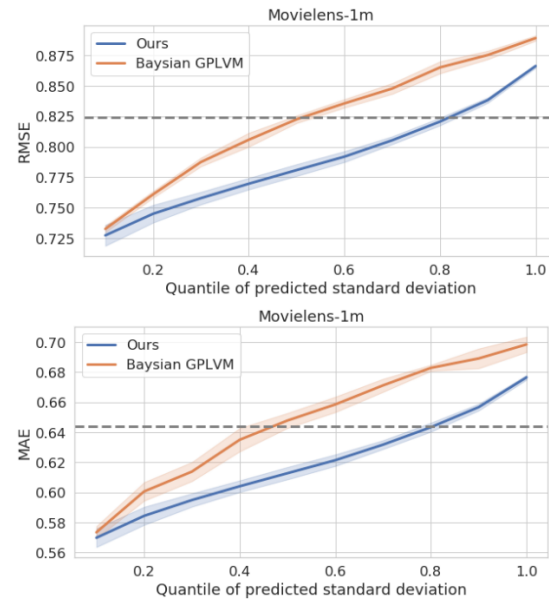


Nickel, Maximilian, et al. A review of relational machine learning for knowledge graphs. *Proceedings of the IEEE* 104.1 (2015): 11-33.

Scalable Gaussian Processes for collaborative filtering

- Our proposal: let f be a Gaussian Process regression and learn the embeddings jointly.
- Challenges and solutions:
 - Kernels \rightarrow Define two kernels for A and B respectively.
 - Scalability (up to 10 million) \rightarrow SOTA scalable GP.
 - Inducing points \rightarrow “*inducing pairs*”, initialized via kernel PCA
 - Evaluation \rightarrow Quantile-Performance-Plot

	RMSE	MAE
CF-NADE [Zheng et al., 2016]	0.829	-
GPLVM [Lawrence and Urtasun, 2009]	0.880	0.644
Sparse FC [Muller et al., 2018]	0.824	-
GC [Berg et al., 2017]	0.832	-
CWOCF [Lu et al., 2013]	0.958	0.761
LLORMA [Lee et al., 2013]	0.865	-
Biased MF	0.863	0.678
SVD++	0.893	0.705
Bayesian GPLVM at $q = 100\%$	0.889	0.698
MW-GP at $q = 100\%$	0.866	0.676
MW-GP at $q = 90\%$	0.838	0.657
MW-GP at $q = 80\%$	0.821	0.643



Yang and Buettner. Multi-output gaussian processes for uncertainty-aware recommender systems. *UAI*. 2021.

Scalable Gaussian Processes for Deep Neural Networks

- CNN + Scalable GP:

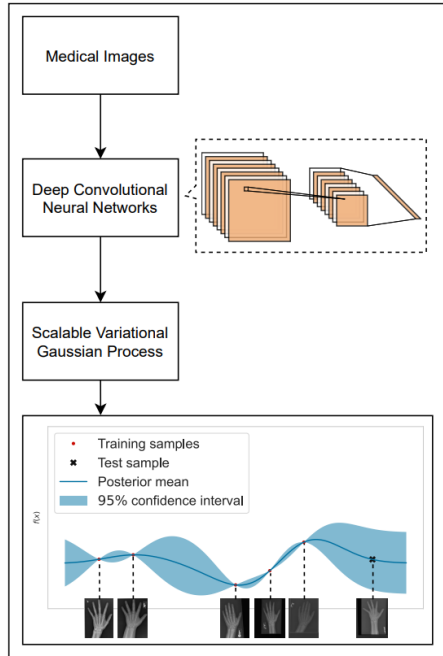
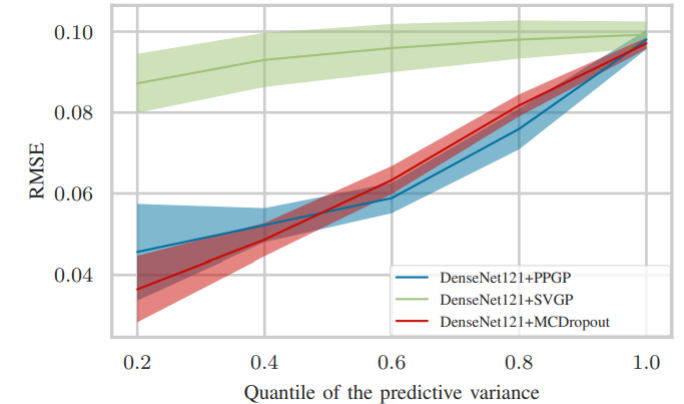
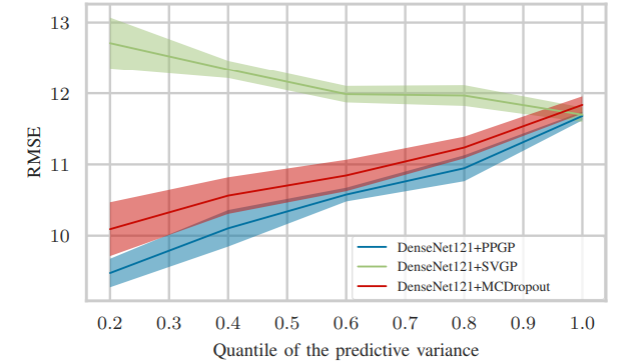


TABLE I
BONE AGE PREDICTION WITH DENSENET121

Output Layer	Transfer Learning	RMSE (No pre-training)	RMSE (DML)	RMSE (CAE)
Linear*	Yes	12.118 ± 0.277	11.667 ± 0.231	14.076 ± 0.281
SVGP†	Yes	11.697 ± 0.102	11.440 ± 0.132	13.536 ± 0.279
PPGP†	Yes	11.679 ± 0.061	11.529 ± 0.089	13.694 ± 0.274
Linear*	No	19.934 ± 0.246	15.805 ± 0.157	15.340 ± 0.390
SVGP†	No	17.723 ± 0.298	15.832 ± 0.284	15.323 ± 0.411
PPGP†	No	18.341 ± 0.234	16.084 ± 0.336	15.752 ± 0.352

TABLE II
LESION LOCALIZATION WITH DENSENET121

Output Layer	Transfer Learning	RMSE (No pre-training)	RMSE (Metric)	RMSE (CAE)
Linear*	Yes	0.102 ± 0.002	0.101 ± 0.002	0.102 ± 0.003
SVGP†	Yes	0.099 ± 0.003	0.101 ± 0.003	0.104 ± 0.004
PPGP†	Yes	0.098 ± 0.002	0.098 ± 0.002	0.099 ± 0.002
Linear*	No	0.116 ± 0.001	0.114 ± 0.003	0.114 ± 0.002
SVGP†	No	0.118 ± 0.002	0.114 ± 0.003	0.112 ± 0.003
PPGP†	No	0.115 ± 0.002	0.111 ± 0.005	0.110 ± 0.002

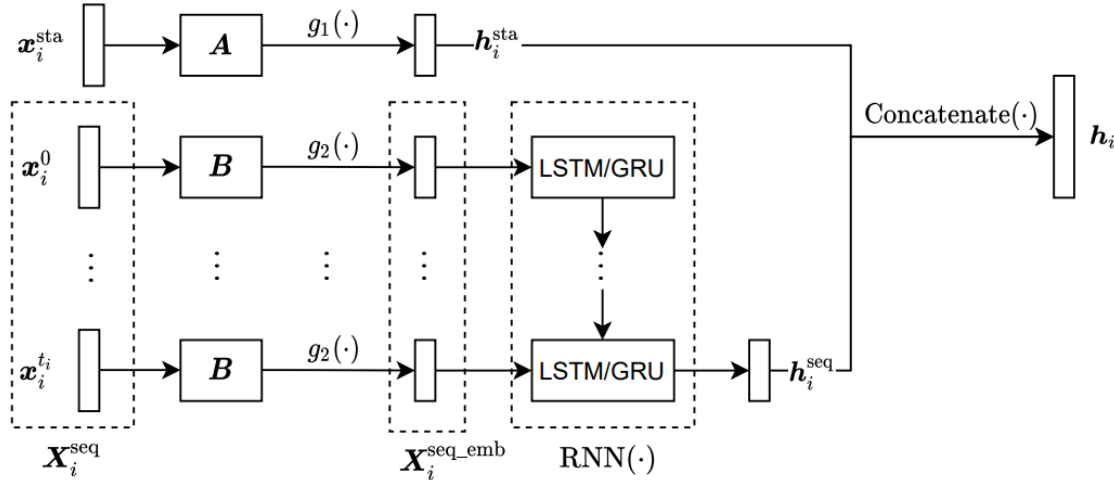


- Tasks: Bone-age prediction, lesion localization.
- Challenge: inducing points initialization.
- Solution: pre-training with auto-encoders.

Wu, et al. Quantifying predictive uncertainty in medical image analysis with deep kernel learning. *2021 IEEE ICHI*. 2021.

Scalable Gaussian Processes for Deep Neural Networks

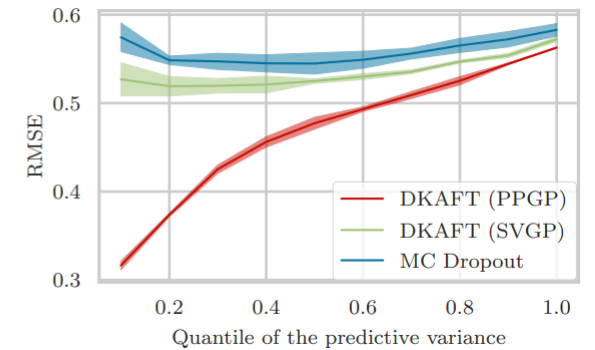
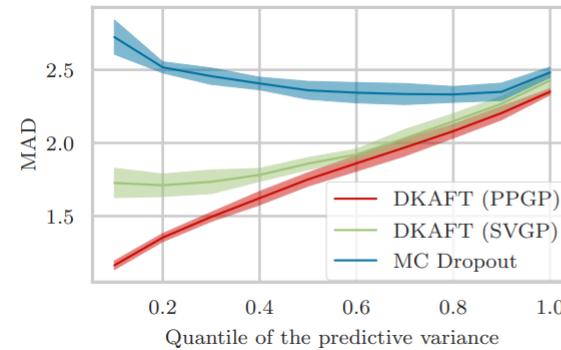
- RNN + Scalable GP



- Tasks: Progression-free survival, length of stay [Sources]
- Challenge: inducing points initialization
- Solution: pre-training, esp. the deep metric learning:

$$\mathcal{J}_{\text{triplet}} = \sum_{i=1}^n [d(h_i^A, h_i^P) - d(h_i^A, h_i^N) + \alpha]_+,$$

Method	Pre-training	Progression-Free Survival		Length-of-Stay	
		MAD	RMSE	MAD	RMSE
Cox Regression	—*	200.800 ± 16.984	1.609 ± 0.054	2.727 ± 0.007	0.638 ± 0.0002
AFT Regression	—*	206.065 ± 8.988	1.685 ± 0.080	2.742 ± 0.014	0.630 ± 0.0001
RNN+AFT	None	150.918 ± 3.009	1.273 ± 0.019	2.476 ± 0.040	0.575 ± 0.003
DKAFT (ExactGP)	None	144.622 ± 8.689	1.225 ± 0.022	—†	—†
DKAFT (SVGP)	None	154.237 ± 13.490	1.211 ± 0.020	2.428 ± 0.056	0.572 ± 0.003
DKAFT (PPGP)	None	147.108 ± 6.284	1.220 ± 0.019	2.351 ± 0.021	0.563 ± 0.001
RNN+AFT	DML	138.155 ± 7.496	1.267 ± 0.007	2.452 ± 0.057	0.568 ± 0.001
DKAFT (ExactGP)	DML	134.422 ± 7.255	1.202 ± 0.012	—†	—†
DKAFT (SVGP)	DML	151.852 ± 11.305	1.221 ± 0.007	2.438 ± 0.079	0.567 ± 0.005
DKAFT (PPGP)	DML	146.616 ± 17.109	1.195 ± 0.008	2.346 ± 0.042	0.557 ± 0.002



Wu, et al. Uncertainty-Aware Time-to-Event Prediction using Deep Kernel Accelerated Failure Time Models. *MLHC (JMLR)*. 2021.

Disclaimer

© Siemens 2023

Subject to changes and errors. The information given in this document only contains general descriptions and/or performance features which may not always specifically reflect those described, or which may undergo modification in the course of further development of the products. The requested performance features are binding only when they are expressly agreed upon in the concluded contract.

All product designations may be trademarks or other rights of Siemens AG, its affiliated companies or other companies whose use by third parties for their own purposes could violate the rights of the respective owner.

| Contact

Published by Siemens 2023

Dr. Yinchong Yang

Senior Key Expert Research Scientist
T DAI HCA-DE

Otto-Hahn-Ring 6
85579 München
Germany