Interpretable representation of latent variables in matrix factorizations

Bachelor's Thesis - Presentation

Lam Ngo

Faculty of Computer Science, University of Vienna

September 30, 2021

Introduction





Figure: Why does Trump appear when one googles "idiot" ? 1

 $^{^{1}} https://www.youtube.com/watch?v=o6zfp6lRw2E\&t=33s\&ab_channel=PBSNew-sHour \\ \\ ^{\square} \wedge ^$

Overview



- 1 Why is interpretability important (in deep learning models)?
- 2 Latent Factor Models and Interpretability
 - Autoencoder
 - Data
 - Interpretable Lens Variable Model
- 3 Experiments
- 4 Conclusion



- Interpretability: relationship to something humanly understandable established²
- Deep learning models come close to the accuracy of human perception but are black box models ³
- GDPR: data subject has right to obtain meaningful information about the logic involved in automated decision making ⁴

²Adel, Ghahramani, and Weller, "Discovering Interpretable Representations for Both Deep Generative and Discriminative Models" (Oct. 2018).

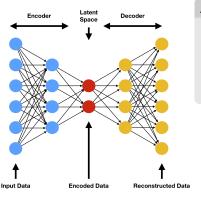
³Chakraborty et al., "Interpretability of deep learning models: A survey of results" (2017).

⁴2018 reform of EU data protection rules (May 25, 2018). ← → ← ≥ → ← ≥ → ← ≥ → ← ≥

Problem - opaque latent factor models



Latent Factor Models and Interpretability Autoencoder University of Vienna



Autoencoder

- $\operatorname{arg\,min}_{A,B}$ $\mathbb{E}[\Delta(x,B\circ A(x))]$
- $A: \mathbb{R}^n \mapsto \mathbb{R}^p \text{ and } B: \mathbb{R}^p \mapsto \mathbb{R}^n$
- Goal: make latent space interpretable

Figure: Architecture of an Autoencoder⁵

- 1 Latent factor model can predict if student answers questions correctly
- 2 Data is from Eedi, an online education platform ⁶
- 3 Students answered diagnostic mathematical questions
- 4 Approx. 15m student-question pairs, 27k questions and 120k students
- Meta data available

Data - Student/Question pairs



Latent Factor Models and Interpretability Data

	QuestionId	Userld	Answerld	IsCorrect	CorrectAnswer	AnswerValue
0	16997	65967	12453206	0	4	2
1	16531	62121	15686710	1	1	1
2	15911	50013	13598796	0	3	1
3	1701	104909	10511925	0	4	3
4	22896	21748	941747	0	1	4

Data - Question Meta Data



Latent Factor Models and Interpretability Data

	QuestionId	SubjectId
0	13090	[3, 32, 71, 77, 141, 185, 186, 214]
1	1855	[3, 71, 75, 86, 178]
2	10423	[3, 32, 38, 239]
3	2290	[3, 32, 33, 144]
4	12785	[3, 32, 33, 144]

Data - Topic Meta Data



Latent Factor Models and Interpretability Data

Su	ubjectld	Name	Parentld	Level
0	3	Maths	NaN	0
1	32	Number	3.0	1
2	33	BIDMAS	144.0	3
3	34	Upper and Lower Bounds	141.0	3
4	35	Calculator Use	32.0	2

- 1 Train the matrix factorization model
 - $\bullet \hat{c}_{ij} = s_i^T q_j + b_i + b_j$
 - lacktriangledown $c_{ij} = ext{question j answered correctly by student i}$
 - \bullet s_i and q_j are latent factors
 - lacksquare b_i and b_j are biases
- 2 Train a relationship between latent space and side information
 - $\hat{c}_{ij} = g(s_i) \times (topics_j * difficulty_j)$
 - g is a linear layer
 - $difficulty_j$ is the share of correct answers for question j of all students

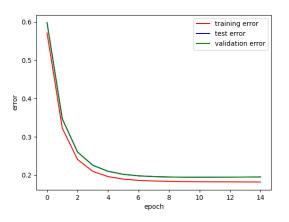


Figure: Error convergence



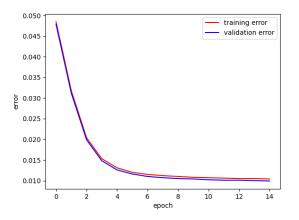


Figure: Error convergence



Experiments University of Vienna

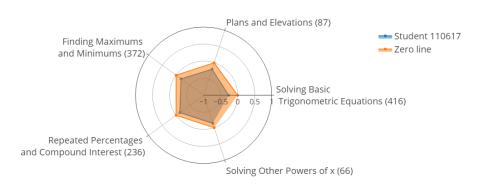


Figure: Lowest knowledge values for Student 110617





Experiments University of Vienna

Subject	Knowledge value	Mean(IsCor- rect)	Mean(diffi- culty) wrong	#questions
			answers	
416	-0.26	0.00	0.37	1
87	-0.20	0.66	0.53	9
372	-0.18	0.00	0.46	1
236	-0.13	1	-	2
66	-0.13	0.5	0.51	4

Table: Lowest knowledge values for Student 110617



Experiments University of Vienna

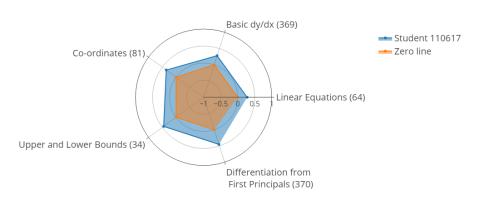


Figure: Highest knowledge values for Student 110617





Experiments University of Vienna

Subject	Knowledge value	Mean(IsCorrect)	Mean(difficulty) wrong answers	#questions
64	0.28	0.75	0.51	8
369	0.28	1	_	2
81	0.36	0.76	0.48	17
34	0.45	0.77	0.47	56
370	0.45	1	-	1

Table: Highest knowledge values for Student 110617

Knowledge Correlations



Experiments University of Vienna

Subject1	Subject2	Corre- lation
Fractions, Decimals and	Fractions	0.88
Percentage Equivalence		
Writing and Simplifying	Perimeter and	0.86
Expressions	Area	
Inequalities on Number	Solving Equations	0.85
Lines		
Indices, Powers and	Volume of Prisms	0.84
Roots		
Formula	Data Representa-	0.83
	tion	

Table: Highest Correlations between Subjects







Conclusion University of Vienna

- Main contribution is a simple way to obtain interpretability in matrix factorization
- Major limitation: no full matrix imputation which would lead to more precise difficulty measure⁷
- 3 Possible extensions: advanced matrix factorization methods, other side information

⁷Wang et al., "Educational Question Mining At Scale: Prediction, Analysis and Personalization" (2021).



References University of Vienna

2018 reform of EU data protection rules.

https://ec.europa.eu/commission/sites/beta-political/files/data-protection-factsheet-changes_en.pdf. European Commission. May 25, 2018. (Visited on 09/17/2021).

- Adel, Tameem, Zoubin Ghahramani, and Adrian Weller. "Discovering Interpretable Representations for Both Deep Generative and Discriminative Models". In: *Proceedings of the 35th International Conference on Machine Learning*. Ed. by Jennifer Dy and Andreas Krause. Vol. 80. Proceedings of Machine Learning Research. PMLR, Oct. 2018, pp. 50–59. URL:
 - http://proceedings.mlr.press/v80/adel18a.html.
- Chakraborty, Supriyo et al. "Interpretability of deep learning models: A survey of results". In: 2017 IEEE SmartWorld, Ubiquitous Intelligence Computing, Advanced Trusted Computed, Scalable Computing Communications, Cloud Big Data Computing, Internet of People and Smart City Innovation



- (SmartWorld/SCALCOM/UIC/ATC/CBDCom/IOP/SCI). 2017, pp. 1–6. DOI: 10.1109/UIC-ATC.2017.8397411.
- Diagnostic Questions The NeurIPS 2020 Education Challenge. https://competitions.codalab.org/competitions/25449. 2020. (Visited on 09/17/2021).
- Wang, Zichao et al. "Educational Question Mining At Scale: Prediction, Analysis and Personalization". In: *Proceedings of the AAAI Conference on Artificial Intelligence*. Vol. 35. 17. 2021, pp. 15669–15677.