

# Knowledge Tracing Machines: Factorization Machines for Knowledge Tracing

Jill-Jênn Vie    Hisashi Kashima



京都大学  
KYOTO UNIVERSITY

KJMLW, February 22, 2019

<https://arxiv.org/abs/1811.03388>

# Practical intro

When exercises are too easy/difficult,  
students get bored/discouraged.

To personalize assessment,

⇒ need a **model** of how people respond to exercises.

## Example

To personalize this presentation,

⇒ need a model of how people respond to my slides.

# Practical intro

When exercises are too easy/difficult,  
students get bored/discouraged.

To personalize assessment,

→ need a **model** of how people respond to exercises.

## Example

To personalize this presentation,

→ need a model of how people respond to my slides.

# Practical intro

When exercises are too easy/difficult,  
students get bored/discouraged.

To personalize assessment,

→ need a **model** of how people respond to exercises.

## Example

To personalize this presentation,

→ need a model of how people respond to my slides.

$p(\text{understanding})$

Practical: 0.9

Theoretical: 0.6

# Theoretical intro

Let us assume  $x$  is **sparse**.

# Theoretical intro

Let us assume  $\mathbf{x}$  is **sparse**.

Linear regression  $y = \langle \mathbf{w}, \mathbf{x} \rangle$

Logistic regression  $y = \sigma(\langle \mathbf{w}, \mathbf{x} \rangle)$  where  $\sigma$  is sigmoid.

Neural network  $x^{(L+1)} = \sigma(\langle \mathbf{w}, \mathbf{x}^{(L)} \rangle)$  where  $\sigma$  is ReLU.

What if  $\sigma : x \mapsto x^2$  for example?

# Theoretical intro

Let us assume  $\mathbf{x}$  is **sparse**.

Linear regression  $y = \langle \mathbf{w}, \mathbf{x} \rangle$

Logistic regression  $y = \sigma(\langle \mathbf{w}, \mathbf{x} \rangle)$  where  $\sigma$  is sigmoid.

Neural network  $x^{(L+1)} = \sigma(\langle \mathbf{w}, \mathbf{x}^{(L)} \rangle)$  where  $\sigma$  is ReLU.

What if  $\sigma : x \mapsto x^2$  for example?

Polynomial kernel  $y = \sigma(1 + \langle \mathbf{w}, \mathbf{x} \rangle)$  where  $\sigma$  is a monomial.

Factorization machine  $y = \langle \mathbf{w}, \mathbf{x} \rangle + \|V\mathbf{x}\|^2$

Mathieu Blondel, Masakazu Ishihata, Akinori Fujino, and Naonori Ueda (2016).  
“Polynomial networks and factorization machines: new insights and efficient  
training algorithms”. In: *Proceedings of the 33rd International Conference on  
International Conference on Machine Learning-Volume 48*. JMLR. org,  
pp. 850–858

# Practical intro

When exercises are too easy/difficult,  
students get bored/discouraged.

To personalize assessment,

→ need a **model** of how people respond to exercises.

## Example

To personalize this presentation,

→ need a model of how people respond to my slides.



# Practical intro

When exercises are too easy/difficult,  
students get bored/discouraged.

To personalize assessment,

→ need a **model** of how people respond to exercises.

## Example

To personalize this presentation,

→ need a model of how people respond to my slides.

$p(\text{understanding})$

Practical: 0.9

Theoretical: 0.9

# Students try exercises

## Math Learning

Items	$5 - 5 = ?$
New student	○

# Students try exercises

## Math Learning

Items	$5 - 5 = ?$	$17 - 3 = ?$
New student	○	○

# Students try exercises

## Math Learning

Items	$5 - 5 = ?$	$17 - 3 = ?$	$13 - 7 = ?$
New student	○	○	×

# Students try exercises

## Math Learning

Items	$5 - 5 = ?$	$17 - 3 = ?$	$13 - 7 = ?$
New student	○	○	×

## Language Learning

	PRON	VERB	PRON	NOUN	CONJ	PRON	VERB	PRON	NOUN
<b>correct:</b>	She	is	my	mother	and	he	is	my	father
<b>student:</b>	she	is		mader	and	he	is		fhader
<b>label:</b>	○	○	×	×	○	○	○	×	×

# Students try exercises

## Math Learning

Items	$5 - 5 = ?$	$17 - 3 = ?$	$13 - 7 = ?$
New student	○	○	✗

## Language Learning

	PRON	VERB	PRON	NOUN	CONJ	PRON	VERB	PRON	NOUN
<b>correct:</b>	She	is	my	mother	and	he	is	my	father
<b>student:</b>	she	is		mader	and	he	is		fhader
<b>label:</b>	○	○	✗	✗	○	○	○	✗	✗

## Challenges

- Users can attempt a same item multiple times
- Users learn over time
- People can make mistakes that do not reflect their knowledge

# Predicting student performance: knowledge tracing

## Data

A population of users answering items

- Events: “User  $i$  answered item  $j$  correctly/incorrectly”

Side information

- If we know the skills required to solve each item e.g., +, ×
- Class ID, school ID, etc.

## Goal: classification problem

Predict the performance of new users on existing items \ Metric:

AUC

## Method

Learn parameters of questions from historical data e.g., *difficulty*

Measure parameters of new students e.g., *expertise*

# Existing work

Model	Basically	Original AUC
Bayesian Knowledge Tracing (Corbett and Anderson 1994)	Hidden Markov Model	0.67



# Existing work

Model	Basically	Original AUC
Bayesian Knowledge Tracing (Corbett and Anderson 1994)	Hidden Markov Model	0.67
Deep Knowledge Tracing (Piech et al. 2015)	Recurrent Neural Network	0.86

# Existing work

Model	Basically	Original AUC	Fixed AUC
Bayesian Knowledge Tracing (Corbett and Anderson 1994)	Hidden Markov Model	0.67	0.63
Deep Knowledge Tracing (Piech et al. 2015)	Recurrent Neural Network	0.86	0.75

# Existing work

Model	Basically	Original AUC	Fixed AUC
Bayesian Knowledge Tracing (Corbett and Anderson 1994)	Hidden Markov Model	0.67	0.63
Deep Knowledge Tracing (Piech et al. 2015)	Recurrent Neural Network	0.86	0.75
Item Response Theory (Rasch 1960) (Wilson et al., 2016)	Online Logistic Regression		0.76

## Existing work

Model	Basically	Original AUC	Fixed AUC
Bayesian Knowledge Tracing (Corbett and Anderson 1994)	Hidden Markov Model	0.67	0.63
Deep Knowledge Tracing (Piech et al. 2015)	Recurrent Neural Network	0.86	0.75
Item Response Theory (Rasch 1960) (Wilson et al., 2016)	Online Logistic Regression		0.76

$$\underbrace{\text{PFA}}_{\text{LogReg}} \leq \underbrace{\text{DKT}}_{\text{LSTM}} \leq \underbrace{\text{IRT}}_{\text{LogReg}} \leq \underbrace{\text{KTM}}_{\text{FM}}$$

# Limitations and contributions

- Several models for knowledge tracing were developed independently
- In our paper, we prove that our approach is more generic

## Our contributions

- Knowledge Tracing Machines unify most existing models
  - Encoding student data to sparse features
  - Then running logistic regression or factorization machines
- Better models found
  - It is better to estimate a bias per item, not only per skill
  - Side information improves performance more than higher dim.

# Our small dataset

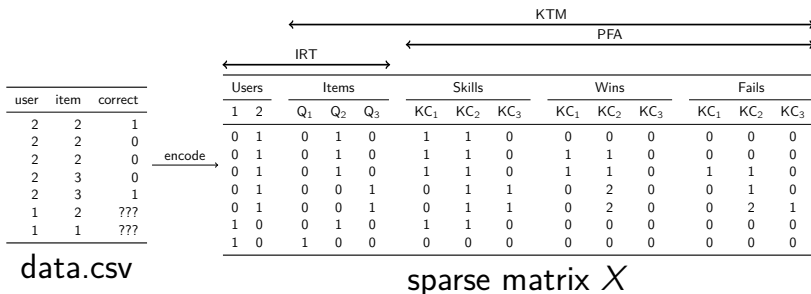
- User 1 answered Item 1 correct
- User 1 answered Item 2 incorrect
- User 2 answered Item 1 incorrect
- User 2 answered Item 1 correct
- User 2 answered Item 2 ???

user	item	correct
1	1	1
1	2	0
2	1	0
2	1	1
2	2	???

dummy.csv

# Our approach

- Encode data to sparse features



- Run logistic regression or factorization machines  
⇒ recover existing models or better models

# Model 1: Item Response Theory

Learn abilities  $\theta_i$  for each user  $i$

Learn easiness  $e_j$  for each item  $j$  such that:

$$Pr(\text{User } i \text{ Item } j \text{ OK}) = \sigma(\theta_i + e_j) \quad \sigma : x \mapsto 1/(1 + \exp(-x))$$

$$\text{logit } Pr(\text{User } i \text{ Item } j \text{ OK}) = \theta_i + e_j$$

Really popular model, used for the PISA assessment

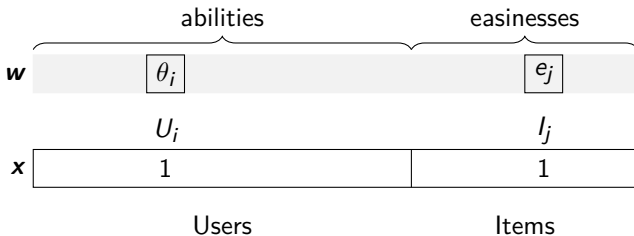
Logistic regression

Learn  $\mathbf{w}$  such that  $\text{logit } Pr(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$



# Graphically: IRT as logistic regression

Encoding “User  $i$  answered Item  $j$ ” with **sparse features**:



$$\langle \mathbf{w}, \mathbf{x} \rangle = \theta_i + e_j = \text{logit } Pr(\text{User } i \text{ Item } j \text{ OK})$$

# Encoding into sparse features

Users			Items		
$U_0$	$U_1$	$U_2$	$I_0$	$I_1$	$I_2$
0	1	0	0	1	0
0	1	0	0	0	1
0	0	1	0	1	0
0	0	1	0	1	0
0	0	1	0	0	1

Then logistic regression can be run on the sparse features.

# Oh, there's a problem

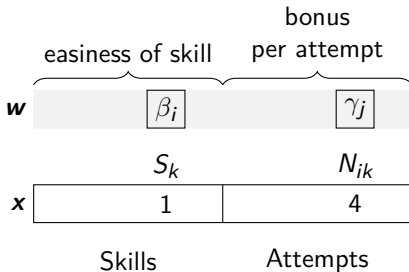
	Users			Items			$y_{\text{pred}}$	$y$
	$U_0$	$U_1$	$U_2$	$I_0$	$I_1$	$I_2$		
User 1 Item 1 OK	0	1	0	0	1	0	0.575135	1
User 1 Item 2 NOK	0	1	0	0	0	1	0.395036	0
User 2 Item 1 <b>NOK</b>	0	0	1	0	1	0	<b>0.545417</b>	<b>0</b>
User 2 Item 1 <b>OK</b>	0	0	1	0	1	0	<b>0.545417</b>	<b>1</b>
User 2 Item 2 NOK	0	0	1	0	0	1	0.366595	0

We predict the same thing when there are several attempts.

# Count number of attempts: AFM

Keep a counter of attempts at skill level:

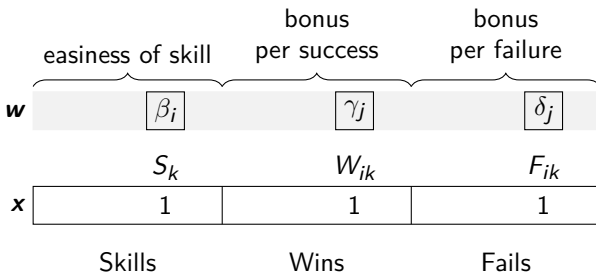
user	item	skill	correct	attempts (for the same skill)
1	1	1	1	0
1	2	2	0	0
2	1	1	0	0
2	1	1	1	1
2	2	2	0	0



# Count successes and failures: PFA

Count separately successes  $W_{ik}$  and fails  $F_{ik}$  of student  $i$  over skill  $k$ .

user	item	skill	correct	wins	fails
1	1	1	1	0	0
1	2	2	0	0	0
2	1	1	0	0	0
2	1	1	1	0	1
2	2	2	0	0	0



## Model 2: Performance Factor Analysis

$W_{ik}$ : how many successes of user  $i$  over skill  $k$  ( $F_{ik}$ : #failures)

Learn  $\beta_k, \gamma_k, \delta_k$  for each skill  $k$  such that:

$$\text{logit } Pr(\text{User } i \text{ Item } j \text{ OK}) = \sum_{\text{Skill } k \text{ of Item } j} \beta_k + W_{ik}\gamma_k + F_{ik}\delta_k$$

Skills			Wins			Fails		
$S_0$	$S_1$	$S_2$	$S_0$	$S_1$	$S_2$	$S_0$	$S_1$	$S_2$
0	1	0	0	0	0	0	0	0
0	0	1	0	0	0	0	0	0
0	1	0	0	0	0	0	0	0
0	1	0	0	0	0	0	1	0
0	0	1	0	0	0	0	0	0

## Better!

	Skills			Wins			Fails			$y_{\text{pred}}$	$y$
	$S_0$	$S_1$	$S_2$	$S_0$	$S_1$	$S_2$	$S_0$	$S_1$	$S_2$		
User 1 Item 1 OK	0	1	0	0	0	0	0	0	0	0.544	1
User 1 Item 2 NOK	0	0	1	0	0	0	0	0	0	0.381	0
User 2 Item 1 <b>NOK</b>	0	1	0	0	0	0	0	0	0	<b>0.544</b>	<b>0</b>
User 2 Item 1 <b>OK</b>	0	1	0	0	0	0	0	1	0	<b>0.633</b>	<b>1</b>
User 2 Item 2 NOK	0	0	1	0	0	0	0	0	0	0.381	0

# Test on a large dataset: Assistments 2009

346860 attempts of 4217 students over 26688 items on 123 skills.

model	dim	AUC	improvement
PFA: skills, wins, fails	0	0.685	+0.07
AFM: skills, attempts	0	0.616	



## Model 3: a new model (but still logistic regression)

model	dim	AUC	improvement
<b>KTM: items, skills, wins, fails</b>	<b>0</b>	<b>0.746</b>	<b>+0.06</b>
IRT: users, items	0	0.691	
PFA: skills, wins, fails	0	0.685	+0.07
AFM: skills, attempts	0	0.616	

# Here comes a new challenger

How to model **pairwise interactions** with **side information**?

## Logistic Regression

Learn a 1-dim **bias** for each feature (each user, item, etc.)

## Factorization Machines

Learn a 1-dim **bias** and a  $k$ -dim **embedding** for each feature

# How to model pairwise interactions with side information?

If you know user  $i$  attempted item  $j$  on **mobile** (not desktop)

How to model it?

$y$ : score of event “user  $i$  solves correctly item  $j$ ”

IRT

$$y = \theta_i + e_j$$

Multidimensional IRT (similar to collaborative filtering)

$$y = \theta_i + e_j + \langle \mathbf{v}_{\text{user } i}, \mathbf{v}_{\text{item } j} \rangle$$

# How to model pairwise interactions with side information?

If you know user  $i$  attempted item  $j$  on **mobile** (not desktop)  
How to model it?

$y$ : score of event “user  $i$  solves correctly item  $j$ ”

IRT

$$y = \theta_i + e_j$$

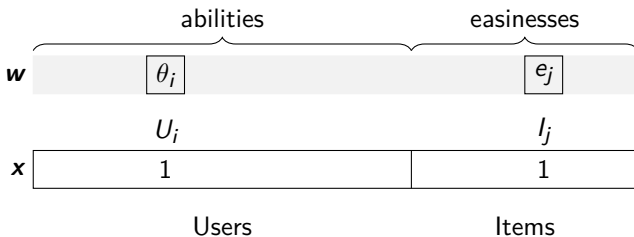
Multidimensional IRT (similar to collaborative filtering)

$$y = \theta_i + e_j + \langle \mathbf{v}_{\text{user } i}, \mathbf{v}_{\text{item } j} \rangle$$

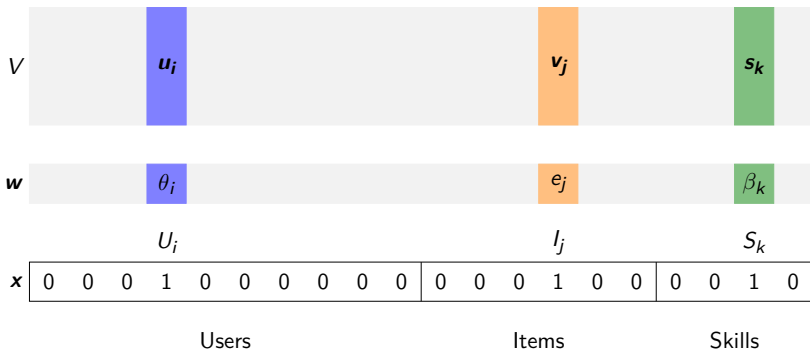
With side information

$$y = \theta_i + e_j + w_{\text{mobile}} + \langle \mathbf{v}_{\text{user } i}, \mathbf{v}_{\text{item } j} \rangle + \langle \mathbf{v}_{\text{user } i}, \mathbf{v}_{\text{mobile}} \rangle + \langle \mathbf{v}_{\text{item } j}, \mathbf{v}_{\text{mobile}} \rangle$$

# Graphically: logistic regression

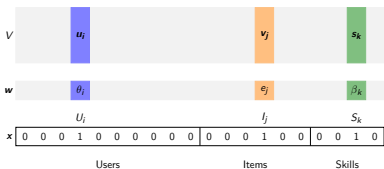


# Graphically: factorization machines



# Formally: factorization machines

Each **user**, **item**, **skill**  $k$  is modeled by bias  $w_k$  and embedding  $v_k$ .



$$\text{logit } p(\mathbf{x}) = \mu + \underbrace{\sum_{k=1}^N w_k x_k}_{\text{logistic regression}} + \underbrace{\sum_{1 \leq k < l \leq N} x_k x_l \langle \mathbf{v}_k, \mathbf{v}_l \rangle}_{\text{pairwise relationships}}$$

Steffen Rendle (2012). “Factorization Machines with libFM”. In: *ACM Transactions on Intelligent Systems and Technology (TIST)* 3.3, 57:1–57:22. DOI: 10.1145/2168752.2168771

# Training using MCMC

Priors:  $w_k \sim \mathcal{N}(\mu_0, 1/\lambda_0)$     $\mathbf{v}_k \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$

Hyperpriors:  $\mu_0, \dots, \mu_n \sim \mathcal{N}(0, 1), \lambda_0, \dots, \lambda_n \sim \Gamma(1, 1) = U(0, 1)$

---

**Algorithm 1** MCMC implementation of FMs

---

**for** each iteration **do**

    Sample hyperp.  $(\lambda_i, \mu_i)_i$  from posterior using Gibbs sampling

    Sample weights  $\mathbf{w}$

    Sample vectors  $\mathbf{V}$

    Sample predictions  $\mathbf{y}$

**end for**

---

Implementation in C++ (libFM) with Python wrapper (pyWFM).

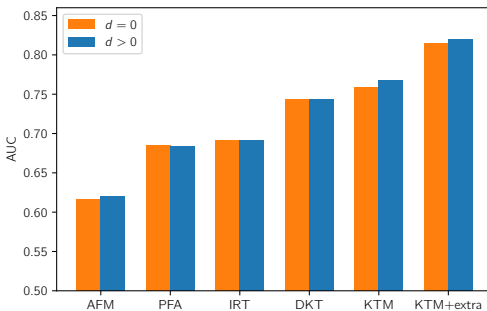
Steffen Rendle (2012). “Factorization Machines with libFM”. In: *ACM Transactions on Intelligent Systems and Technology (TIST)* 3.3, 57:1–57:22. DOI: 10.1145/2168752.2168771



# Datasets

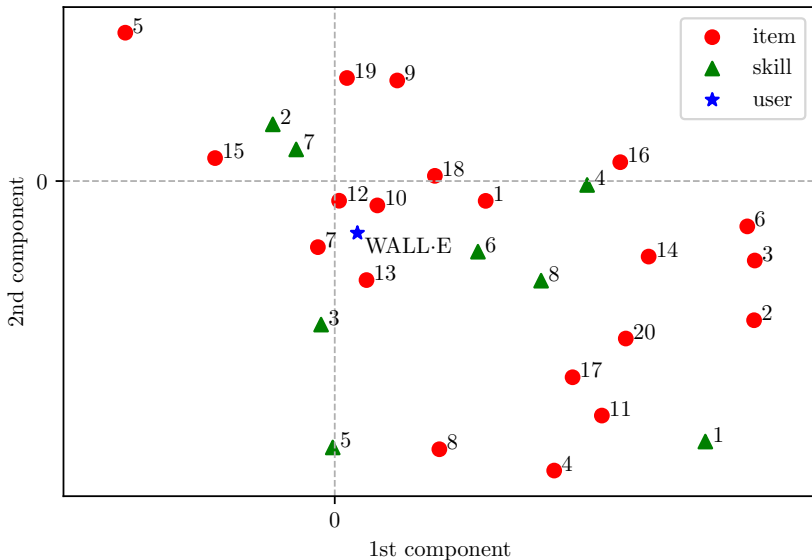
Name	Users	Items	Skills	Skills/i	Entries	Sparsity	Attempts/u
fraction	536	20	8	2.800	10720	0.000	1.000
timss	757	23	13	1.652	17411	0.000	1.000
ecpe	2922	28	3	1.321	81816	0.000	1.000
assistentments	4217	26688	123	0.796	346860	0.997	1.014
berkeley	1730	234	29	1.000	562201	0.269	1.901
castor	58939	17	2	1.471	1001963	0.000	1.000

# AUC results on the Assistments dataset



model	dim	AUC	improvement
KTM: items, skills, wins, fails, extra	5	0.819	
KTM: items, skills, wins, fails, extra	0	0.815	+0.05
KTM: items, skills, wins, fails	10	0.767	
KTM: items, skills, wins, fails	0	0.759	+0.02
<i>DKT</i> (Wilson et al., 2016)	100	0.743	+0.05
IRT: users, items	0	0.691	
PFA: skills, wins, fails	0	0.685	+0.07
AFM: skills, attempts	0	0.616	

# Bonus: interpreting the learned embeddings

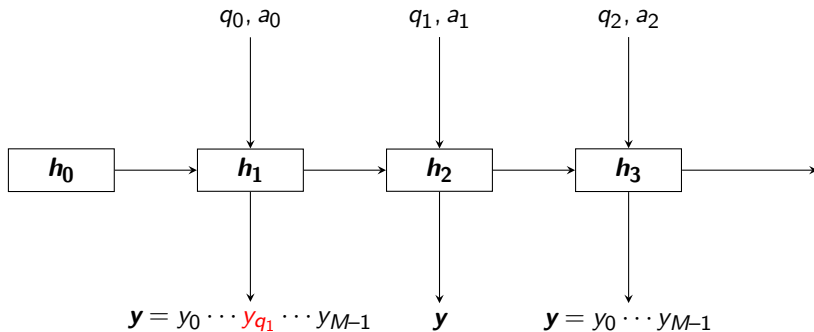


# What 'bout recurrent neural networks?

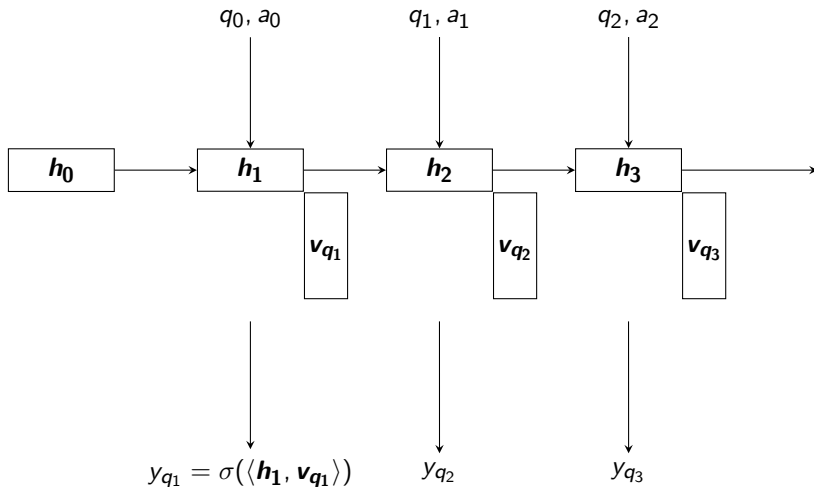
Deep Knowledge Tracing: model the problem as sequence prediction

- Each student on skill  $q_t$  has performance  $a_t$
- How to predict outcomes  $\mathbf{y}$  on every skill  $k$ ?
- Spoiler: by measuring the evolution of a latent state  $h_t$

# Graphically: deep knowledge tracing



# Graphically: there is a MIRT in my DKT



## Drawback of Deep Knowledge Tracing

DKT does not model individual differences.

Actually, Wilson even managed to beat DKT with (1-dim!) IRT.

By estimating on-the-fly the student's learning ability, we managed to get a better model.

AUC	BKT	IRT	PFA	DKT	DKT-DSC
Assistments 2009	0.67	0.75	0.70	0.73	0.91
Assistments 2012	0.61	0.74	0.67	0.72	0.87
Assistments 2014	0.64	0.67	0.69	0.72	0.87
Cognitive Tutor	0.61	0.81	0.76	0.79	0.81

Sein Minn, Yi Yu, Michel Desmarais, Feida Zhu, and Jill-Jênn Vie (2018). "Deep Knowledge Tracing and Dynamic Student Classification for Knowledge Tracing". In: *Proceedings of the 18th IEEE International Conference on Data Mining*, to appear. URL: <https://arxiv.org/abs/1809.08713>

# Take home message

**Knowledge tracing machines** unify many existing EDM models

- Side information improves performance more than higher  $d$
- We can visualize learning (and provide feedback to learners)

Already provides better results than vanilla **deep neural networks**

- Can be combined with FMs



# Do you have any questions?

Read our article:

Knowledge Tracing Machines

<https://arxiv.org/abs/1811.03388>

Try our tutorial:

<https://github.com/jilljenn/ktm>

I'm interested in:

- predicting student performance
- recommender systems
- optimizing human learning using reinforcement learning

[vie@jill-jenn.net](mailto:vie@jill-jenn.net)



Blondel, Mathieu, Masakazu Ishihata, Akinori Fujino, and Naonori Ueda (2016). “Polynomial networks and factorization machines: new insights and efficient training algorithms”. In: *Proceedings of the 33rd International Conference on International Conference on Machine Learning-Volume 48*. JMLR. org, pp. 850–858.



Corbett, Albert T and John R Anderson (1994). “Knowledge tracing: Modeling the acquisition of procedural knowledge”. In: *User modeling and user-adapted interaction 4.4*, pp. 253–278.



Minn, Sein, Yi Yu, Michel Desmarais, Feida Zhu, and Jill-Jênn Vie (2018). “Deep Knowledge Tracing and Dynamic Student Classification for Knowledge Tracing”. In: *Proceedings of the 18th IEEE International Conference on Data Mining*, to appear. URL: <https://arxiv.org/abs/1809.08713>.



Piech, Chris, Jonathan Bassen, Jonathan Huang, Surya Ganguli, Mehran Sahami, Leonidas J Guibas, and Jascha Sohl-Dickstein (2015). “Deep knowledge tracing”. In: *Advances in Neural Information Processing Systems (NIPS)*, pp. 505–513.



Rasch, Georg (1960). “Studies in mathematical psychology: I. Probabilistic models for some intelligence and attainment tests.”. In:



Rendle, Steffen (2012). “Factorization Machines with libFM”. In: *ACM Transactions on Intelligent Systems and Technology (TIST)* 3.3, 57:1–57:22. DOI: 10.1145/2168752.2168771.