

DAS3H: Modeling Student Learning and Forgetting for Optimally Scheduling Distributed Practice of Skills

Benoît Choffin, Fabrice Popineau, Yolaine Bourda & Jill-Jênn Vie

LRI/CentraleSupélec - University of Paris-Saclay | RIKEN AIP, now Inria Lille

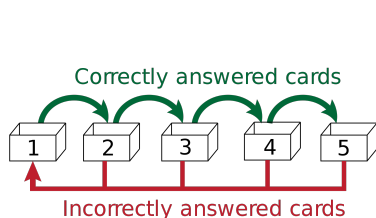


Inria Bordeaux | November 27, 2019

Mitigating human forgetting with spaced repetition

- Human learners face a constant trade-off between **acquiring new knowledge** and **reviewing old knowledge**
- Cognitive science provides simple + robust learning strategies for improving LT memory
 - Spaced repetition
 - Testing
- Can we do better? **Yes**, by providing students with an *adaptive* and *personalized* spacing scheduler.

Mitigating human forgetting with spaced repetition



memorizing

暗記

あんき



Model-based

Ex. select the item whose memory strength is closest to a threshold θ [Lindsey, Shroyer, Pashler, and Mozer 2014] → “almost forgotten”

Model-free

Bandit methods such as [Clement, Roy, Oudeyer, and Lopes 2013] (of course), other reinforcement learning methods

Beyond flashcard memorization

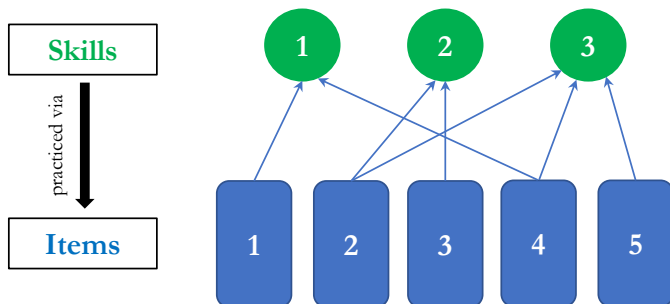
Problem: these algorithms are designed for optimizing *pure memorization* (of facts, vocabulary, . . .)

- In real-world educational settings, students also need to learn to master and remember a set of **skills**
- In that case, specific items are the only way to practice one or multiple skills because *we do not have to memorize the content directly*
- Traditional adaptive spacing schedulers are **not applicable for learning skills**

Extension to skill practice and review

Item-skill relationships require expert labor and are synthesized inside a binary q-matrix →

	skill 1	skill 2	skill 3
item 1	1	0	0
item 2	0	1	1
item 3	0	1	0
item 4	1	0	1
item 5	0	0	1



Limitations of student models

We need to be able to infer skill memory strength and dynamics, however in the student modeling literature:

- some models leverage item-skills relationships
- some others incorporate forgetting

But none does both!

Our contribution

We take a model-based approach for this task.

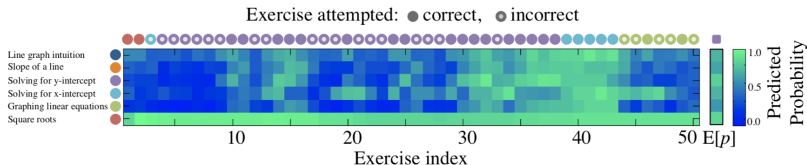
- ① Traditional adaptive spacing algorithms can be extended to review and practice skills (not only flashcards).
- ② We developed a new student *learning* and *forgetting* model that leverages item-skill relationships: **DAS3H**.
 - DAS3H outperforms 4 SOTA student models on 3 datasets.
 - Incorporating skill info + forgetting effect improves over models that consider one or the other.
 - Using precise temporal information on past skill practice + assuming different learning/forgetting curves **for different skills** improves performance.

Outline

- ① Knowledge tracing
- ② Our model DAS3H
- ③ Experiments
- ④ Conclusion

Knowledge tracing

Predict future student performance given their history

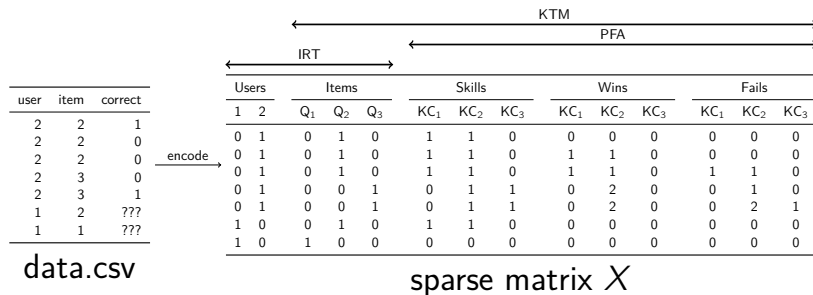


Given $(q_t, a_t)_{t \leq T}$ for former students

q_t is the question index, $a_t \in \{0, 1\}$ is the correctness

For new students, given $(q_t, a_t)_{t \leq T}$ and q_{T+1} , guess a_{T+1}

Encoding data into sparse features



Then run logistic regression or factorization machines (teaser)

Model 1: Item Response Theory

Learn abilities θ_i for each user i

Learn easiness e_j for each item j such that:

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

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

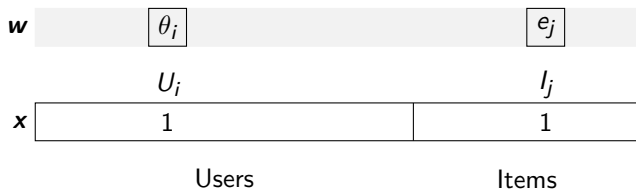
Logistic regression

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

Usually with L2 regularization: $\|\mathbf{w}\|_2^2$ penalty \leftrightarrow Gaussian prior

Graphically: IRT as logistic regression

Encoding of “User i answered Item j ” into 2-hot vectors



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

Encoding

```
python encode.py --users --items
```

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

data/dummy/X-ui.npz

Then logistic regression can be run on the sparse features:

```
python lr.py data/dummy/X-ui.npz
```

Oh, there's a problem

```
python encode.py --users --items
```

```
python lr.py data/dummy/X-ui.npz
```

	Users			Items			y_{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 NOK	0	0	1	0	1	0	0.545417	0
User 2 Item 1 OK	0	0	1	0	1	0	0.545417	1
User 2 Item 2 NOK	0	0	1	0	0	1	0.366595	0

We predict the same thing when there are several attempts.

⇒ Need temporal features

Count successes and failures

Keep track of what the student has done before:

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

data/dummy/data.csv

Model 2: Performance Factor Analysis

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

Learn β_k , γ_k , δ_k for each skill k such that:

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

`python encode.py --skills --wins --fails`

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

`data/dummy/X-swf.npz`

Better!

```
python encode.py --skills --wins --fails
```

```
python lr.py data/dummy/X-swf.npz
```

			Skills			Wins			Fails			y_{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	NOK	0	1	0	0	0	0	0	0	0	0.544	0
User 2	Item 1	OK	0	1	0	0	0	0	0	1	0	0.633	1
User 2	Item 2	NOK	0	0	1	0	0	0	0	0	0	0.381	0

Model 3: DASH

→ DASH = item **D**ifficulty, student **A**bility, and **S**tudent **H**istory

DASH [Lindsey, Shroyer, Pashler, and Mozer 2014] bridges the gap between *Factor Analysis models* and *memory models*:

$$\mathbb{P}(Y_{s,j,t} = 1) = \sigma(\alpha_s - \delta_j + h_{\theta}(t_{s,j,1:\ell}, y_{s,j,1:\ell-1}))$$

where:

- $Y_{s,j,t}$ binary correctness of student s answering item j at time t ;
- σ logistic function;
- α_s ability of student s ;
- δ_j difficulty of item j ;
- h_{θ} summarizes the effect of the $\ell - 1$ previous attempts of s on j at times $t_{s,j,1:\ell-1}$ + the binary outcomes $y_{s,j,1:\ell-1}$.

DASH

Lindsey et al. chose:

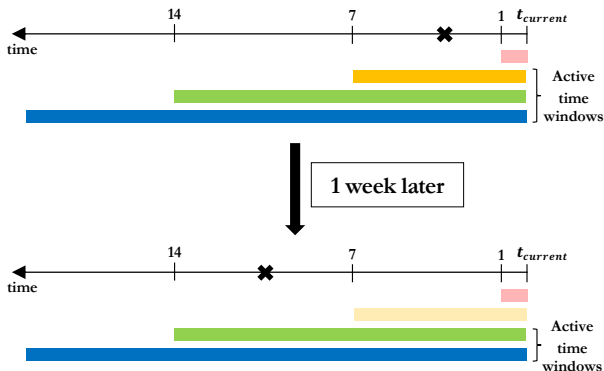
$$h_{\theta}(t_{s,j,1:l}, y_{s,j,1:l-1}) = \sum_{w=0}^{W-1} \theta_{2w+1} \log(1 + c_{s,j,w}) - \theta_{2w+2} \log(1 + a_{s,j,w})$$

where:

- w indexes a set of expanding **time windows**;
- $c_{s,j,w}$ number of correct answers of s on j in time window w ;
- $a_{s,j,w}$ number of attempts of s on j in time window w ;
- θ is *learned* by DASH.

DASH

Assuming that the set of time windows is $\{1, 7, 14, +\infty\}$:



DASH

DASH:

- accounts for both *learning* and *forgetting* processes;
- induces diminishing returns of practice inside a time window (log-counts);
- has a time module h_θ inspired by ACT-R [Anderson, Matessa, and Lebiere 1997] and MCM [Pashler, Cepeda, Lindsey, Vul, and Mozer 2009].

From DASH to DAS3H

- DASH
 - outperforms a hierarchical Bayesian IRT on Lindsey et al. experimental data (vocabulary learning).
 - was successfully used to adaptively personalize item review in a real-world cognitive psychology experiment.
- However, DASH
 - does not handle multiple skill item tagging → useful to account for knowledge transfer from one item to another.
 - assumes that memory decays at the same rate for every KC.

Our model DAS3H

We extend DASH in **3 ways**:

- ① Extension to handle multiple skills tagging: new temporal module h_θ that also takes the multiple skills into account.
 - Influence of the temporal distribution of past attempts and outcomes can differ from one skill to another.
- ② Estimation of easiness parameters for *each* item j and skill k ;
- ③ Use of KTMs [Vie and Kashima 2019] instead of mere logistic regression for multidimensional feature embeddings and pairwise interactions.

Our model DAS3H

→ DAS3H = item **D**ifficulty, student **A**bility, **S**kill and **S**tudent **S**kill practice **H**istory

For an embedding dimension of $d = 0$, DAS3H is:

$$\mathbb{P}(Y_{s,j,t} = 1) = \sigma(\alpha_s - \delta_j + \underbrace{\sum_{k \in KC(j)} \beta_k}_{\text{skill easiness biases}} + h_\theta(t_{s,j,1:l}, y_{s,j,1:l-1})).$$

We choose:

$$h_\theta(t_{s,j,1:l}, y_{s,j,1:l-1}) = \sum_{k \in KC(j)} \sum_{w=0}^{W-1} \theta_{k,2w+1} \log(1 + c_{s,k,w}) - \theta_{k,2w+2} \log(1 + a_{s,k,w}).$$

→ Now, h_θ can be seen as a sum of *skill* memory strengths!

Learning multidimensional feature embeddings

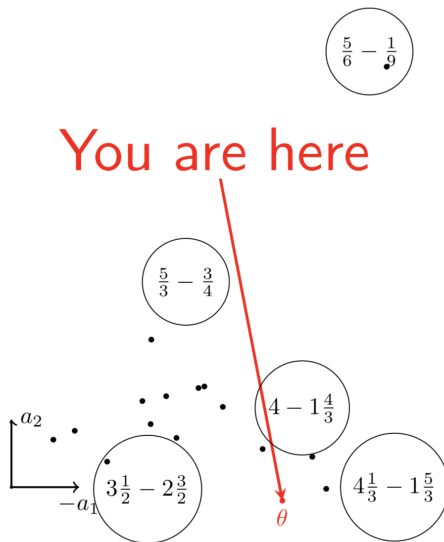
Logistic Regression

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

Factorization Machines

Learn a **bias** and an **embedding** for each feature

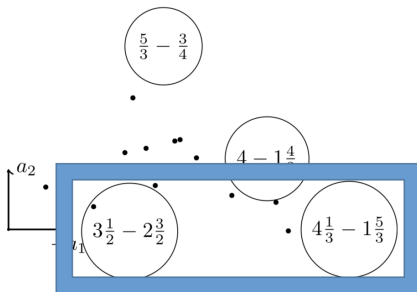
What can be done with multidimensional embeddings?



Interpreting the components

$$\frac{5}{6} - \frac{1}{9}$$

**Items that
discriminate
only over one dimension**



$$3\frac{1}{2} - 2\frac{3}{2}$$

$$b = 0.13$$
$$-a_1 = 2.01$$
$$a_2 = -0.03$$

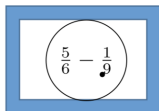
$$4\frac{1}{3} - 2\frac{4}{3}$$

$$b = -0.46$$
$$-a_1 = 4.65$$
$$a_2 = -0.02$$

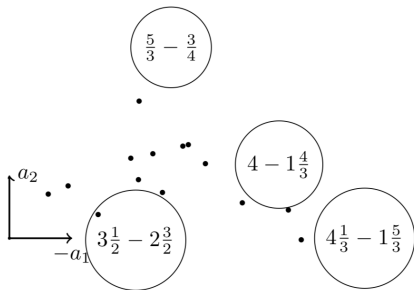
$$4\frac{1}{3} - 1\frac{5}{3}$$

$$b = -1.99$$
$$-a_1 = 5.66$$
$$a_2 = 0.00$$

Interpreting the components



**Items that
highly discriminate
over both dimensions**



$$\frac{3}{4} - \frac{3}{8}$$

$$b = 1.09$$
$$-a_1 = 5.54$$
$$a_2 = 6.22$$

$$\frac{5}{6} - \frac{1}{9}$$

$$b = -0.28$$
$$-a_1 = 5.29$$
$$a_2 = 6.44$$

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)

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

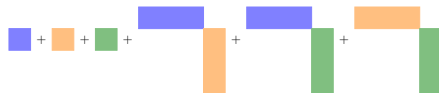
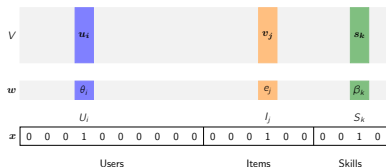
With side information

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

Knowledge Tracing Machines (KTMs)

Just pick features (ex. **user**, **item**, **skill**) and you get a student model

Each feature k is modeled by bias w_k and embedding v_k .



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

Jill-Jênn Vie and Hisashi Kashima (2019). “Knowledge Tracing Machines: Factorization Machines for Knowledge Tracing”. In: *Proceedings of the 33th AAAI Conference on Artificial Intelligence*, to appear. URL: <http://arxiv.org/abs/1811.03388>

Experiments

- ① Experimental setting
- ② Contenders
- ③ Datasets
- ④ Main results
- ⑤ Further analyses

Experimental setting

How to compare ML models?

Train the models on one part of the dataset

Test on the other part

Gather prediction metrics, compare the models

- **5-fold cross-validation** at the student level: predicting binary outcomes on **unseen** students (*strong generalization*)
- Distributional assumptions to **avoid overfitting**:
 - When $d = 0$: L2 regularization/ $\mathcal{N}(0, 1)$ prior
 - When $d > 0$: hierarchical distributional scheme
- Same time windows as Lindsey et al.: $\{1/24, 1, 7, 30, +\infty\}$

Contenders

5 contenders:

- **DAS3H**
- DASH [Lindsey, Shroyer, Pashler, and Mozer 2014]
- IRT/MIRT [Linden and Hambleton 2013]
- PFA [Pavlik, Cen, and K. R. Koedinger 2009]
- AFM [Cen, K. Koedinger, and Junker 2006]

Every model was cast within the KTM framework → 3 embedding dimensions (0, 5 & 20) + sparse feature encoding.

	users	items	skills	wins	fails	attempts	tw [KC]	tw [items]
DAS3H	x	x	x	x		x	x	
DASH	x	x		x		x		x
IRT/MIRT	x	x						
PFA			x	x	x			
AFM			x			x		

Datasets

- 3 datasets: ASSISTments 2012-2013, Bridge to Algebra 2006-2007 & Algebra I 2005-2006 (KDD Cup 2010)
 - Data consists of logs of student-item interactions on 2 ITS
 - Selected because they contain *both* timestamps and items with multiple skills → rare species in the EDM datasets fauna
- Preprocessing scheme: removed users with < 10 interactions, interactions with NaN skills, duplicates

Dataset	Users	Items	Skills	Interactions	Mean correctness	Skills per item	Mean skill delay	Mean study period
assist12	24,750	52,976	265	2,692,889	0.696	1.000	8.54	98.3
bridge06	1,135	129,263	493	1,817,427	0.832	1.013	0.83	149.5
algebra05	569	173,113	112	607,000	0.755	1.363	3.36	109.9

Table 2: Datasets characteristics

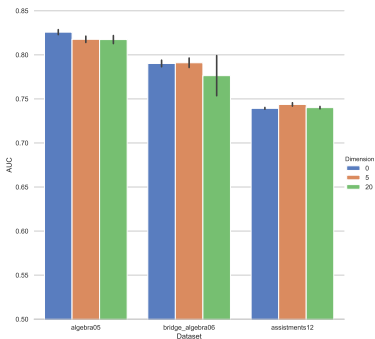
Main results

model	algebra05	bridge06	assist12
DAS3H	0.826 \pm 0.003	0.790 \pm 0.004	0.739 \pm 0.001
DASH	0.773 \pm 0.002	0.749 \pm 0.002	0.703 \pm 0.002
IRT	0.771 \pm 0.007	0.747 \pm 0.002	0.702 \pm 0.001
PFA	0.744 \pm 0.004	0.739 \pm 0.003	0.668 \pm 0.002
AFM	0.707 \pm 0.005	0.692 \pm 0.002	0.608 \pm 0.002

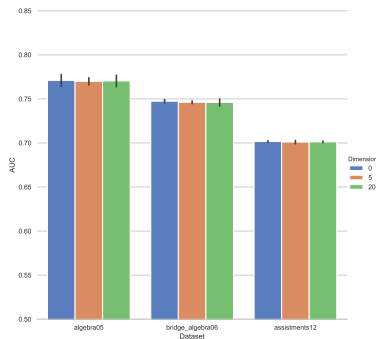
Table 3: AUC comparison between the different student models for an embedding dimension $d = 0$ (all datasets, 5-fold cross-validation).

→ On every dataset, **DAS3H outperforms** the other models (between +0.04 and +0.05 AUC compared to DASH).

Main results



(a) DAS3H



(b) IRT

Figure 1: AUC comparison on two models for $d = 0, 5$ and 20 (all datasets, 5-fold cross-validation).

→ The impact of the multidim feature embeddings is small and not consistent across datasets and models (+ unstable sometimes).

Importance of time windows

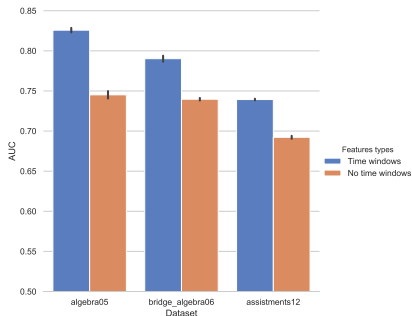


Figure 2: AUC comparison on DAS3H *with* and *without* time windows features (all datasets, 5-fold cross-validation).

Without time windows, h_θ counts past wins and attempts in DAS3H.

→ Using **temporal distribution of past skill practice** instead of simple win/fail counters improves AUC performance: the **when** matters.

Importance of different learning/forgetting curves per skill

	d	bridge06	algebra05	assist12
DAS3H	0	0.790 \pm 0.004	0.826 \pm 0.003	0.739 \pm 0.001
	5	0.791 \pm 0.005	0.818 \pm 0.004	0.744 \pm 0.002
	20	0.776 \pm 0.023	0.817 \pm 0.005	0.740 \pm 0.001
DAS3H _{1p}	0	0.757 \pm 0.003	0.789 \pm 0.009	0.701 \pm 0.002
	5	0.757 \pm 0.005	0.787 \pm 0.005	0.700 \pm 0.001
	20	0.757 \pm 0.003	0.789 \pm 0.006	0.701 (<1e-3)

Table 4: AUC comparison between DAS3H and DAS3H_{1p} (all datasets, 5-fold cross-validation).

→ Assuming **different learning and forgetting curves for different skills** in DAS3H consistently yields better predictive power: some skills are easier to learn and slower to forget.

In a nutshell

- Human forgetting is *ubiquitous* but luckily:
 - **Cognitive science** gives us efficient and simple learning strategies
 - **ML** can build us tools to **personalize these strategies** and further improve LT memory retention
- Adaptive spacing algorithms have been focusing on *pure memorization* (e.g. vocabulary learning)
 - They can be used for **optimizing practice and retention of skills**
- Our student model **DAS3H**
 - incorporates information on *skills* **and** *forgetting* to predict learner performance
 - shows higher predictive power than other SOTA student models
 - fits our model-based approach for optimally scheduling skill review

Thanks for your attention!

Our paper is already available at:

<https://arxiv.org/abs/1905.06873>

Python code is freely available on our GitHub pages:

<https://github.com/BenoitChoffin/das3h>

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

To send us questions:

benoit.choffin@lri.fr

jill-jenn.vie@inria.fr



Anderson, John R, Michael Matessa, and Christian Lebiere (1997). “ACT-R: A theory of higher level cognition and its relation to visual attention”. In: *Human-Computer Interaction* 12.4, pp. 439–462.



Cen, Hao, Kenneth Koedinger, and Brian Junker (2006). “Learning factors analysis—a general method for cognitive model evaluation and improvement”. In: *International Conference on Intelligent Tutoring Systems*. Springer, pp. 164–175.



Clement, Benjamin, Didier Roy, Pierre-Yves Oudeyer, and Manuel Lopes (2013). “Multi-armed bandits for intelligent tutoring systems”. In: *arXiv preprint arXiv:1310.3174*.



Linden, Wim J van der and Ronald K Hambleton (2013). *Handbook of modern item response theory*. Springer Science & Business Media.



Lindsey, Robert V, Jeffery D Shroyer, Harold Pashler, and Michael C Mozer (2014). “Improving students’ long-term knowledge retention through personalized review”. In: *Psychological science* 25.3, pp. 639–647.



Pashler, Harold, Nicholas Cepeda, Robert V Lindsey, Ed Vul, and Michael C Mozer (2009). “Predicting the optimal spacing of study: A multiscale context model of memory”. In: *Advances in neural information processing systems*, pp. 1321–1329.



Pavlik, Philip I., Hao Cen, and Kenneth R. Koedinger (2009). “Performance Factors Analysis - A New Alternative to Knowledge Tracing”. In: *Proceedings of the 14th International Conference on Artificial Intelligence in Education, AIED 2009*, pp. 531–538.



Vie, Jill-Jênn and Hisashi Kashima (2019). “Knowledge Tracing Machines: Factorization Machines for Knowledge Tracing”. In: *Proceedings of the 33th AAAI Conference on Artificial Intelligence*, to appear. URL: <http://arxiv.org/abs/1811.03388>.