Distinguishing Humans and GenAl Using Multiple-Choice Tests (and IRT)

Remember at March 2023...



How a tiny company with few rules is making fake images go mainstream Midjourney, the year-old firm behind recent fake visuals of Trump and the pope, illustrates the lack of oversight

Generative AI Era Challenge



In Education

סכנולוגיה ומדע

בגלל הבינה המלאכותית: המוסדות האקדמיים בבעיה רצינית

מחקר חדש בבריטניה חושף: %94 מהעבודות האקדמיות שנכתבו על ידי בינה מלאכותית לא זוהו על ידי המרצים • ״זו בעיה רצינית שמפחיתה מערכם של תעודות ותארים״

> מערכת ישראל היום סייראל 2/12/2024, פי:09, עודכן 2/12/2024, פי:09

> > חדשות • חינוך וחברה

עבודת סמינר בלחיצת כפתור: הבינה המלאכותית מתחילה לשנות את האקדמיה

שמירה 🚺 קריאת זן 📃



Compromising

- Learning
- Assessment

In Education

Open questions/essay



Multiple choice questions (MCQ's)

Unexplored

Allegedly "limited bits information"

Classically the view is that **the whole test** measures something with noise.



In Education

Open questions/essay



Multiple choice questions (MCQ's)

Unexplored

Allegedly "limited bits information"



Classically the view is th**ot she whole test** measures something with noise.

Approach- use Item Response Theory (IRT)

• Test\instrument - a collection of Items that measure the hidden trait(s)



Response matrix:



Response matrix:



Response matrix:



For more students and items we get this kind of response matrix:

Response matrix:

	Items a b c d e f g h i j k l 1 1 1 0 0 0 1 1 0 0 1 1 1 1 0 0 0 0 1 1 0 0 1 1 0 1 0 0 0 0 1 1 0 0 0 1 1 1 0 0 0 0 1 1 1 1 1 1 1 1 0 0 0 0 1 1 1 1 1 1 1 1 0 0 1												
Students	а	b	С	d	е	f	g	h	i	j	k	Т	score
А	1	1	1	0	0	0	0	1	1	0	0	1	6
В	1	0	1	1	0	0	0	0	1	0	0	0	4
С	1	1	1	0	1	0	0	1	1	1	1	1	9
D	1	0	1	1	0	0	0	0	1	0	0	1	5
Е	0	1	1	0	0	1	0	1	1	1	1	1	8
F	1	1	1	1	0	0	0	1	1	0	0	1	7
G	1	0	1	0	0	1	0	0	1	0	1	1	6
Н	1	0	1	0	0	0	0	0	0	0	0	1	3
I	1	1	1	1	0	0	0	1	1	0	0	1	7

We rearrange the matrix:

						lte	ms							Step 1:
Students	а	b	С	d	е	f	g	h	i	j	k	Ι	score	score: high to lo
А	1	1	1	0	0	0	0	1	1	0	0	1	6	
В	1	0	1	1	0	0	0	0	1	0	0	0	4	
С	1	1	1	0	1	0	0	1	1	1	1	1	9	
D	1	0	1	1	0	0	0	0	1	0	0	1	5	
E	0	1	1	0	0	1	0	1	1	1	1	1	8	
F	1	1	1	1	0	0	0	1	1	0	0	1	7	
G	1	0	1	0	0	1	0	0	1	0	1	1	6	
Н	1	0	1	0	0	0	0	0	0	0	0	1	3	
I	1	1	1	1	0	0	0	1	1	0	0	1	7	↓

W

We rearrange the matrix:

						lte	ms							
Students	s a	b	С	d	е	f	g	h	i	j	k	Ι	score	Ability
С	1	1	1	0	1	0	0	1	1	1	1	1	9	
E	0	1	1	0	0	1	0	1	1	1	1	1	8	
F	1	1	1	1	0	0	0	1	1	0	0	1	7	
l I	1	1	1	1	0	0	0	1	1	0	0	1	7	
А	1	1	1	0	0	0	0	1	1	0	0	1	6	
G	1	0	1	0	0	1	0	0	1	0	1	1	6	
D	1	0	1	1	0	0	0	0	1	0	0	1	5	
В	1	0	1	1	0	0	0	0	1	0	0	0	4	
Н	1	0	1	0	0	0	0	0	0	0	0	1	3	¥

We rearrange the matrix:

matrix:	St It	t ep em	2: s: ea	asy	to h	nard						-	W frc	hat is a om the o	n easy item data?
							lte	ms							
Students		а	b	С	d	е	f	g	h	i	j	k	I	score	Ability
С		1	1	1	0	1	0	0	1	1	1	1	1	9	
Е		0	1	1	0	0	1	0	1	1	1	1	1	8	
F		1	1	1	1	0	0	0	1	1	0	0	1	7	
l I		1	1	1	1	0	0	0	1	1	0	0	1	7	
А		1	1	1	0	0	0	0	1	1	0	0	1	6	
G		1	0	1	0	0	1	0	0	1	0	1	1	6	
D		1	0	1	1	0	0	0	0	1	0	0	1	5	
В		1	0	1	1	0	0	0	0	1	0	0	0	4	
Н		1	0	1	0	0	0	0	0	0	0	0	1	3	↓

All students answered correctly

We rearrange the matrix:

natrix:	Ste Iten	p 2: ns: (eas	y to	har	d					→		
						lte	ms						
Students	С	i	а	Ι	е	f	g	h	b	j	k	d	Ability
С	1	1	1	1	1	0	0	1	1	1	1	0	9
Е	1	1	0	1	0	1	0	1	1	1	1	0	8
F	1	1	1	1	0	0	0	1	1	0	0	1	7
I	1	1	1	1	0	0	0	1	1	0	0	1	7
А	1	1	1	1	0	0	0	1	1	0	0	0	6
G	1	1	1	1	0	1	0	0	0	0	1	0	6
D	1	1	1	1	0	0	0	0	0	0	0	1	5
В	1	1	1	0	0	0	0	0	0	0	0	1	4
Н	1	0	1	1	0	0	0	0	0	0	0	0	3

We rearrange the matrix:

matriv	Step	o 2:											
	lten	าร: (easy	y to	har	d							
						lte	ms						
Students	С	i	а	Π	b	h	d	k	f	j	е	g	Ability
С	1	1	1	1	1	1	0	1	0	1	1	0	9
Е	1	1	0	1	1	1	0	1	1	1	0	0	8
F	1	1	1	1	1	1	1	0	0	0	0	0	7
I	1	1	1	1	1	1	1	0	0	0	0	0	7
А	1	1	1	1	1	1	0	0	0	0	0	0	6
G	1	1	1	1	0	0	0	1	1	0	0	0	6
D	1	1	1	1	0	0	1	0	0	0	0	0	5
В	1	1	1	0	0	0	1	0	0	0	0	0	4
Н	1	0	1	1	0	0	0	0	0	0	0	0	3

C+--- 7.

We rearrange the matrix:

natrix:	S It	cer :en	ס ב: רא:	easy	y to	har	d					→		
							lte	ms						
Students		С	i	а	Ι	b	h	d	k	f	j	е	g	Ability
С		1	1	1	1	1	1	0	1	0	1	1	0	9
Е		1	1	0	1	1	1	0	1	1	1	0	0	8
F		1	1	1	1	1	1	1	0	0	0	0	0	7
I		1	1	1	1	1	1	1	0	0	0	0	0	7
А		1	1	1	1	1	1	0	0	0	0	0	0	6
G		1	1	1	1	0	0	0	1	1	0	0	0	6
D		1	1	1	1	0	0	1	0	0	0	0	0	5
В		1	1	1	0	0	0	1	0	0	0	0	0	4
Н		1	0	1	1	0	0	0	0	0	0	0	0	3

Dificulty $\propto \frac{1}{\text{Facility}}$

 Facility
 9
 8
 8
 5
 5
 4
 3
 2
 2
 1
 0



Dificulty
$$\propto \frac{1}{\text{Facility}}$$



What happens with large deviations?



Dificulty
$$\propto \frac{1}{\text{Facility}}$$

What happens with large deviations?



Calculate an anomaly measure



Calculate an anomaly measure



						lte	ms							
Students	С	i	а	Ι	b	h	d	k	f	j	е	g	Ability	Anomaly
С	1	1	1	1	1	1	0	1	0	1	1	0	9	5
E	1	1	0	1	1	1	0	1	1	1	8	0	8	9
F	1	1	0	1	1	1	1	0	0	0	0	0	7	4
I	1	1	1	1	1	1	1	0	0	0	0	0	7	0
А	1	1	1	1	1	1	0	0	0	0	0	0	6	0
G	1	1	1	1	0	0	0	1	1	0	0	0	6	6
D	1	1	1	1	1	1	1	0	0	0	0	0	5	0
В	1	1	1	1	0	0	1	0	0	0	0	0	4	2
Н	1	1	1	1	0	0	0	0	0	0	0	0	3	0
J														
К														
L														
Μ														



						lte	ms							
Students	С	i	а	I	b	h	d	k	f	j	е	g	Ability	Anomaly
С	1	1	1	1	1	1	0	1	0	1	1	0	9	5
Е	1	1	0	1	1	1	0	1	1	1	8	0	8	9
F	1	1	0	1	1	1	1	0	0	0	0	0	7	4
I	1	1	1	1	1	1	1	0	0	0	0	0	7	0
А	1	1	1	1	1	1	0	0	0	0	0	0	6	0
G	1	1	1	1	0	0	0	1	1	0	0	0	6	6
D	1	1	1	1	1	1	1	0	0	0	0	0	5	0
В	1	1	1	1	0	0	1	0	0	0	0	0	4	2
Н	1	1	1	1	0	0	0	0	0	0	0	0	3	0
J														
К														
L														NOT SIN
М														

						lte	ms							
Students	С	i	а	Ι	b	h	d	k	f	j	е	g	Ability	Anomaly
С	1	1	1	1	1	1	0	1	0	1	1	0	9	5
Е	1	1	0	1	1	1	0	1	1	1	8	0	8	9
F	1	1	0	1	1	1	1	0	0	0	0	0	7	4
I	1	1	1	1	1	1	1	0	0	0	0	0	7	0
А	1	1	1	1	1	1	0	0	0	0	0	0	6	0
G	1	1	1	1	0	0	0	1	1	0	0	0	6	6
D	1	1	1	1	1	1	1	0	0	0	0	0	5	0
В	1	1	1	1	0	0	1	0	0	0	0	0	4	2
Н	1	1	1	1	0	0	0	0	0	0	0	0	3	0
J	1	1	0	0	1	0	1	0	1	0	0	1	6	
К	0	0	1	1	0	1	0	0	1	0	1	0	5	
L	0	1	0	0	1	0	1	0	0	1	0	1	5	NOT SIN
Μ	1	0	0	0	1	1	0	0	1	0	0	1	5	



"Guttman error"



Our working assumption:

Chatbots "behave" as students with different latent traits

- Two different tests with students' responses per item
 - Formative Chemistry MCQ online "preparation to matriculation" test
 - Summative Quantitative Psychometrics part*
- Pollution level -5%
- Three premium version chatbots
- Four anomaly measures *G*, *G**,*U*, *ZU3 Person Fit Statistics*

Results

Chatbots "behave" as students with different latent traits



What if the pollution level is increased?

- Results until now: pollution level 5%
- We also tried with 10%, 25%.
- Pollution level $\uparrow \Rightarrow$ Separation humans chats \downarrow

• Widespread use of chatbots *means* harder separation

Do different chatbots behave differently?



By using the same measures we show: Chats can be separated using MCQ test

Who from who?

Depends on the test and on the specific measure

To conclude...

Humans vs. chatbots – •

Have different latent traits; can be distinguished using MCQ

G*

0.5

0.6

0.7

- Chatbots vs. other chatbots ٠
- Widespread use of chatbots means harder separation •



THANKS FOR LISTENING

Link to the paper

