

Leveraging a Cognitive Model to Measure Subjective Similarity of Human and GPT-4 Written Content

Anonymous EMNLP submission

Abstract

Cosine similarity between two documents can be computed using token embeddings formed by Large Language Models (LLMs) such as GPT-4, and used to categorize those documents across a range of uses. However, these similarities are ultimately dependent on the corpora used to train these LLMs, and may not reflect subjective similarity of individuals or how their biases and constraints impact similarity metrics. This lack of cognitively-aware personalization of similarity metrics can be particularly problematic in educational and recommendation settings where there is a limited number of individual judgements of category or preference, and biases can be particularly relevant. To address this, we rely on an integration of an Instance-Based Learning (IBL) cognitive model with LLM embeddings to develop the Instance-Based Individualized Similarity (IBIS) metric. This similarity metric is beneficial in that it takes into account individual biases and constraints in a manner that is grounded in the cognitive mechanisms of decision making. To evaluate the IBIS metric, we also introduce a dataset of human categorizations of emails as being either dangerous (phishing) or safe (ham). This dataset is used to demonstrate the benefits of leveraging a cognitive model to measure the subjective similarity of human participants in an educational setting.

1 Introduction

When humans categorize textual information, such as when giving recommendations or learning to categorize documents, we often use our personal subjective concepts to complete the task. One example of this is giving a recommendation of a funny book to a friend, which requires not only our own subjective conceptualization of humor, but also an understanding of the similarities and differences between ourselves and our friends. While humans perform this task with relative ease, recommendation systems (Ansari et al., 2000) and educational

tools (Nafea et al., 2019) typically do not have personalized measurements of subjective concepts, either for themselves or the people that are using these systems, potentially hindering their efficacy.

The specific use case we are interested in is a learning setting where students are categorizing documents and receiving feedback of the accuracy of their categorization. In this work, we focus specifically on students categorizing emails as being safe (ham) or dangerous (phishing) in a training setting to help users identify and defend against phishing email attacks.

When these systems do incorporate data from human judgements to determine the subjective similarity of, they typically do so by pooling together as many judgements from different people as they can, and aggregate their measurement (Xia et al., 2015). This can be effective from a machine learning perspective, since more data can mean improved performance for the general public. But in terms of providing an individualized experience to students in educational settings or end-users in recommendation settings, this type of data aggregation approach leaves something to be desired.

When methods do attempt to account for individual measures of similarity, they typically employ machine learning based methods (Shojaei and Saneifar, 2021). While these approaches can be beneficial in some use cases, they are not grounded by the biases and constraints inherent in human learning in a way that is afforded through cognitive modeling.

In this work, we propose a method for providing personalized metrics of subjective concepts that can determine the similarity between sets of text, with additional applications in textual categorization and educational feedback. This is done by leveraging a cognitive model of human learning and decision making that can act as a digital twin to individuals, and predict their behavior and opinions on a wider set of stimuli. This cognitive model

043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081
082
083

incorporates LLM embeddings into its prediction of human behavior, allowing for flexible and efficient connections between cognitive models and LLMs.

Other than this method, the main contribution of this work is in presenting a dataset of human participant judgements in an educational setting learning to correctly categorize emails as being either safe ‘ham’ or dangerous ‘phishing’. This is a nuanced and difficult task that could potentially be based on multiple different subjective classifications such as the level of urgency in a text, whether it has a suspicious tone, or whether it is making an offer that seems too good to be true. All of these features are relevant to determining if an email is genuine, and the way that individuals perceive an email as having these features can be highly individualistic.

Alongside data from these human judgements we present a dataset of emails written by human cybersecurity experts, as well as emails generated by GPT-4 while relying on various levels of information from human prompt engineers. The final part of this dataset is a set of conversations between human participants and a GPT-4o model providing feedback to students. These conversations and the resulting educational improvement of students can provide useful insight into the prompting of LLMs for educational settings.

In total this dataset represents 39230 human judgements from 430 participants making decisions while observing a set from 1440 GPT-4 or human generated emails, as well as 20487 messages between human participants and the GPT-4o teacher model.

2 Phishing Email Categorization Dataset

One of the contributions of this work is the presentation of a dataset of human judgements categorizing emails as being either phishing or ham. These emails were created from various sources, including human cybersecurity experts, GPT-4 generation, as well as humans working with GPT-4 in collaboration. Data from human participants categorizing these emails as being either ham or phishing in an educational setting was made available online at ¹. This dataset includes 39230 categorization judgements from 384 human participants of 1440 possible emails.

The second component of this dataset is the emails shown to participants, which were either

written by human cybersecurity experts, a GPT-4 model working alone, or a combination of human and GPT-4 model work. 360 base emails written by human experts were used to form three additional versions of these base emails. These alternative versions included a ‘human-written gpt4-styled’ version that used the email body written by human experts, the ‘gpt4-written and gpt4-styled’ version that was fully rewritten by GPT-4, and the ‘gpt4-written plaintext-styled’ version that stripped the HTML and CSS styling applied by the GPT-4 model. These emails as well as the original prompts to generate them are included in the presented dataset.

The final component of this dataset is a set of conversations between human participants and a GPT-4o model prompted to teach the participant to identify phishing emails. In this experiment three out of the eight experimental conditions involved human participants discussing the emails that they were categorizing with a GPT-4o model. This model was prompted to serve as an educational tool and varied in the type of information that was included in these prompts across experimental conditions. These teacher-student conversations consist of 20487 messages sent between human participants and the GPT-4o model.

3 Background: Cognitive Model

The cognitive model used in this work to predict the subjective similarity of human participants decisions on unseen emails relies on Instance Based Learning Theory (IBLT) (Gonzalez et al., 2003). This learning theory describes the mathematical foundation of cognitive mechanisms that underlie human decision making in dynamic environments, such as learning tasks. Cognitive models that rely on the mathematical framework of IBLT are called Instance-Based Learning models, which are used to define the features relevant for decision making tasks, and to predict the mechanisms of dynamic decision making based on these features.

One of the benefits of employing IBL models over alternatives is that they take into account the past experiences of participants and the impact of limitations like memory size and decay that can bias decision making. IBL models have been applied onto predicting human behavior in dynamic decision making tasks, including repeated binary choice tasks (Gonzalez and Dutt, 2011; Lejarraga et al., 2012), theory of mind applications (Nguyen

¹<https://osf.io/wbg3r/>

and Gonzalez, 2022), and practical applications such as identifying phishing emails (Cranford et al., 2019; Malloy and Gonzalez, 2024), cyber defense (Cranford et al., 2020), and cyber attack decision-making (Aggarwal et al., 2022). The following sections outline the mathematical foundation of IBL models, and gives attention to the method of integrating these concepts into predictions of subjective similarity of categories.

3.1 Activation

IBL models work by storing instances i in memory \mathcal{M} , composed of utility outcomes u_i and options k composed of features j in the set of features \mathcal{F} of environmental decision alternatives. These options are observed in an order represented by the time step t , and the time step that an instance occurred in is given $\mathcal{T}(i)$. IBL models predict the value of options in decision-making tasks by selecting the action that maximizes the value function. In calculating this activation, the similarity between instances in memory and the current instance is represented by summing over all attributes the value S_{ij} , which is the similarity of attribute j of instance i to the current state. This gives the activation equation as:

$$A_i(t) = \ln \left(\sum_{t' \in \mathcal{T}_i(t)} (t - t')^{-d} \right) + \mu \sum_{j \in \mathcal{F}} \omega_j (S_{ij} - 1) + \sigma \xi \quad (1)$$

The parameters that are set either by modelers or set to default values are the decay parameter d ; the mismatch penalty μ ; the attribute weight of each j feature ω_j ; and the noise parameter σ . The default values for these parameters are $(d, \mu, \omega_j, \sigma) = (0.5, 1, 1, 0.25)$. The value ξ is drawn from a normal distribution $\mathcal{N}(-1, 1)$ and multiplied by the noise parameter σ to add random noise to the activation.

3.2 Probability of Retrieval

The probability of retrieval represents the probability that a single instance in memory will be retrieved when estimating the value associated with an option. To calculate this probability of retrieval, IBL models apply a weighted soft-max function onto the memory instance activation values $A_i(t)$ giving the equation:

$$P_i(t) = \frac{\exp A_i(t)/\tau}{\sum_{i' \in \mathcal{M}_k} \exp A_{i'}(t)/\tau} \quad (2)$$

The parameter that is either set by modelers or set to its default value is the temperature parameter τ , which controls the uniformity of the probability distribution defined by this soft-max equation. The default value for this parameter is $\tau = \sigma\sqrt{2}$.

3.3 Blended Value

The blended value of an option k is calculated at time step t according to the utility outcomes u_i weighted by the probability of retrieval of that instance P_i and summing over all instances in memory \mathcal{M}_k to give the equation:

$$V_k(t) = \sum_{i \in \mathcal{M}_k} P_i(t) u_i \quad (3)$$

The blended value of different options is key to predicting the subjective similarity of human participants in a way that is both individualized, and takes into account the experience of students in educational settings. This is done in our proposed individualized metric of subjective similarity by comparing the blended value of an individual that is engaged in a task to categorize emails as being either phishing or spam.

4 Methods of Measuring Similarity

4.1 Human Similarity Measures

We can use the accuracy of human participant categorizations and the confidence that participants selected to their judgement to plot for each email their level of similarity to phishing and ham emails. For both of these metrics, a higher value signifies that participants were more likely to categorize an emails as being a member of that group with a high confidence. These results are graphed in Figure 1, which plots the phishing and ham similarity of each email based on the average of human performance.

The reaction time and confidence weighted subjective similarity of an email x is given by multiplying the probability of a human participant categorizing that email as category c giving $cs(x|c) = p(c|x)r(c|x)c(c|x)$. where $p(c|x)$ is the probability of categorization, $r(c|x)$ is the reaction time normalized to between 0 and 1, and $c(c|x)$ is the confidence additionally normalized to between 0 and 1. This gives the subjective similarity as:

$$HS(x, x') = \frac{cs(x|c)cs(x'|c)}{\sum_{c' \in C} cs(x|c) \sum_{c' \in C} cs(x'|c')} \quad (4)$$

This metric of subjective similarity depends on primarily the categorization of emails from human

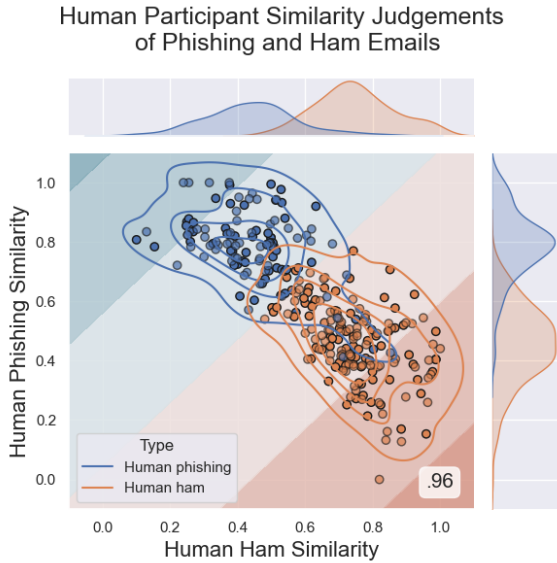


Figure 1: Average Human Similarity measure for phishing (blue) and ham (orange) emails based on the categorization, confidence, and response latency of participant responses. Shaded region represents a logistic regression classifier trained on 100 train-test splits of size 50% with the accuracy shown in the lower right.

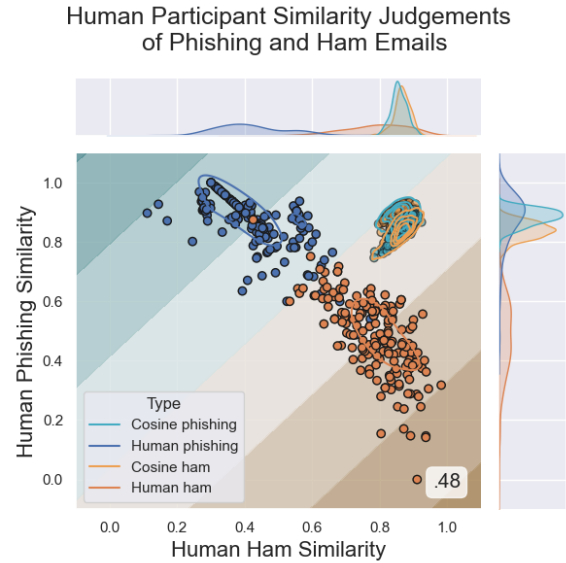


Figure 2: Average cosine and human participant similarity for phishing (light blue) and ham (light orange) emails. Shaded region represents a logistic regression classifier trained on 100 train-test splits of size 50% with the accuracy shown in the lower right.

272 participants, but additionally takes into account
 273 their confidence, with higher confidences of phishing
 274 categorizations indicating an email is more
 275 similar to the phishing email category. Additionally,
 276 this metric takes into account the reaction time
 277 of human participants in making their judgement,
 278 with faster judgements additionally indicating that
 279 a category is more similar to members of that
 280 category average. The goal of the IBIS method is to
 281 reflect this type of subjective similarity, and is compared
 282 to several alternative measures of similarity
 283 described in the following sections.

284 4.2 Cosine Similarity

285 Cosine similarity is the most commonly used met-
 286 ric of similarity of word and document embed-
 287 dings, with many applications from classification
 288 (Park et al., 2020), recommendation systems (Khat-
 289 ter et al., 2021), educational tutorial systems (Wu
 290 et al., 2023), question answering (Aithal et al.,
 291 2021), and more (Patil et al., 2023). However, there
 292 are limitations to using cosine similarity such as
 293 in documents with high-frequency words (Zhou
 294 et al., 2022), and the presence of false information
 295 (Borges et al., 2019), both of which are concerns
 296 for phishing email education.

297 A simple way to apply cosine similarity onto
 298 the task of categorizing an email as being either

299 phishing or ham is to collect a large number of
 300 labelled emails and compute the average of the
 301 embeddings of these labelled emails. Once this
 302 embedding average is collected, we can measure
 303 the cosine distance of any given email embedding
 304 and the average of both categories.

$$\begin{aligned}
 \text{CS}(x, x') &= \frac{x^T x'}{\|x\| \|x'\|} \\
 &= \frac{x^T x'}{\sqrt{x^T x} \sqrt{x'^T x'}}
 \end{aligned}
 \tag{5}$$

305 The cosine similarity of each email embedding
 306 to the mean embedding of that category is shown
 307 in Figure , and compared to our metric of subjective
 308 similarity that is dependent on human partici-
 309 pant categorization, confidence, and reaction time.
 310 From this, we can see that on average the embed-
 311 dings are calculated as being significantly more
 312 similar to each other compared to the subjective
 313 similarities of human participants.
 314

315 Comparing the accuracy of using the cosine simi-
 316 larity metrics, we can see that the logistic regres-
 317 sion of predicting the human subjective similarity
 318 has now decreased in accuracy to 0.48, from the
 319 previous accuracy of 0.96 when forming a logistic
 320 regression of human participant similarity judge-
 321 ments alone. This significant decrease is due to the
 322 gap between cosine similarity and the subjective

Human Participant and Weighted Cosine Similarity of Phishing and Ham Emails

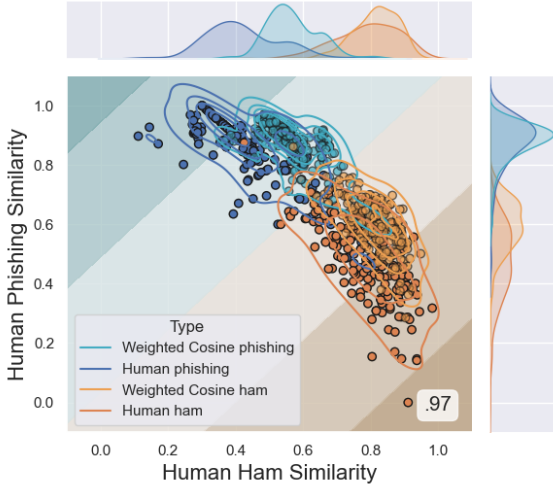


Figure 3: Average weighted cosine and human participant similarity for phishing (light blue) and ham (light orange) emails. Shaded region represents a logistic regression classifier trained on 100 train-test splits of size 50% with the accuracy shown in the lower right.

Human Participant and Pruned Cosine Similarity of Phishing and Ham Emails

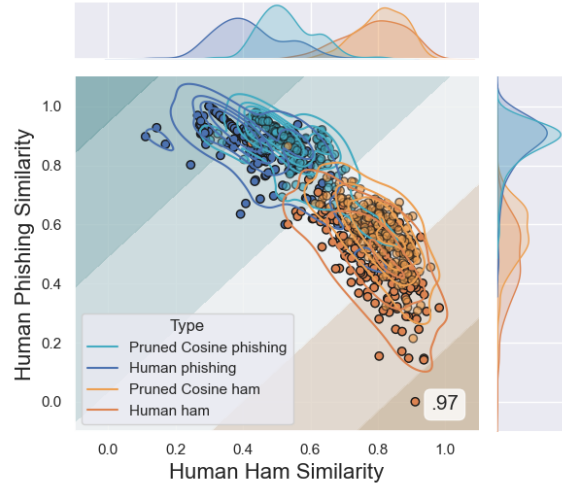


Figure 4: Average pruned cosine and human participant similarity for phishing (light blue) and ham (light orange) emails. Shaded region represents a logistic regression classifier trained on 100 train-test splits of size 50% with the accuracy shown in the lower right.

323 similarity of human participants. One solution to
 324 this is to differentially weight the indices of embed-
 325 ding values, which is explored next.

326 4.3 Weighted Cosine Similarity

327 Distance weighted cosine similarity is a common
 328 method employed in utilizing embeddings (Li and
 329 Han, 2013), which has been applied onto measur-
 330 ing similarity of online instruction in educational
 331 settings (Lahitani et al., 2016), as well as several
 332 cybersecurity specific applications like ransomware
 333 detection (Moussaileb et al., 2021), and inside
 334 attacker detection (Khan et al., 2019). In this
 335 work, we employ weighted cosine similarities of
 336 embeddings formed from emails categorized as
 337 being either ham or phishing, and compare it to
 338 human subjective similarity judgements. This can
 339 be done by defining the weighted cosine similarity
 340 of an email embedding as:

$$\begin{aligned}
 CS_w(x, x', W) &= \frac{(Wx)^T(Wx')}{\|Wx\| \|Wx'\|} \\
 &= \frac{x^T W^T W x'}{\sqrt{x^T W^T W x} \sqrt{x'^T W^T W x'}}
 \end{aligned}
 \tag{6}$$

342 From this definition of the weighted cosine sim-
 343 ilarity, it is relatively simple to construct the em-
 344 bedding weight matrix A in a way that minimizes

345 the mean squared error of the distance between
 346 weighted cosine embeddings and the subjective
 347 similarity metrics of participants. This allows for
 348 the classification of emails in a way that reflects
 349 the confidence and categorization of human partici-
 350 pants in an educational setting. The results of this
 351 weighting are shown in Figure 3, which compares
 352 the average human participant subjective similarity
 353 and the weighted cosine similarity of email embed-
 354 dings.

355 The accuracy of the logistic regression fit to
 356 weighted cosine similarities of phishing and ham
 357 emails when predicting human subjective similarity
 358 has increased to 0.97 from the unweighted accu-
 359 racy of 0.48. These improved similarity metrics
 360 indicate that weighting cosine similarity based on
 361 data from a large dataset of human participants can
 362 result in a metric that more accurately reflects the
 363 average of human subjects' subjective similarity
 364 metrics.

365 4.4 Pruning Document Embeddings

366 The final method of comparison for developing
 367 individualized metrics of similarity is embedding
 368 pruning, where embeddings are reduced in size
 369 based on feedback from human categorizations to
 370 better account for their subjective similarity (Man-
 371 rique et al., 2023). This method was originally de-
 372 signed for word embeddings with a larger number

of categories that are more varied than our application. We adjusted this approach to apply it onto reflecting human categorization of emails into only two related categories of phishing and ham emails.

After making these adjustments to the embedding pruning method the result is a similarity metric calculated by ranking embedding value by how well it predicts the different human categorization performance, and selecting only the top 500 embedding values, representing just under 20% of the size of the embedding, as was done in (Manrique et al., 2023). These top predictive embedding values are retained, while all other values are masked to 0. After this, cosine similarity can be calculated with the standard approach, resulting in the similarity shown in Figure 5.

5 Instance-Based Individualized Similarity (IBIS)

To determine an individual participant’s metric of similarity, we employ an IBL model that is serving as a digital twin of the participant. The result in an Instance-Based Individualized Similarity (IBIS) metric. The benefits of IBIS are in the ability to predict human judgements on unseen documents or feedback from recommendations, and enhance measurements of subjective similarity. Importantly, these predictions of human behavior are not merely relying on a separate machine learning based technique, but rather a cognitive model that is inspired by the human cognitive mechanisms underlying decision making and thus able to account for natural biases and constraints in humans.

$$IBIS(x, x') = \frac{V_k(c|x)V_k(c|x')}{\sum_{c' \in C} V_k(c'|x) \sum_{c' \in C} V_k(c'|x')} \quad (7)$$

Predictions of Instance-Bases Individual Similarity are done using an IBL model that is currently serving as a digital twin with the same experience as an individual participant. Using this we determine the value that the IBL model assigns to predicting a category c as $V_k(c|x)$, or the value the IBL model assigns to choosing option c as the category of document x . Then, we can divide this value by the same categorization value assigned to each alternative categorization of the same document. This results in the IBIS metric which can be calculated after each decision is made by a participant, as shown in the pseduo-code in Algorithm 1.

Input: default utility u_0 , a memory dictionary $\mathcal{M} = \{\}$, global counter $t = 1$, step limit L . Dataset of stimuli D

```

repeat
  Initialize a counter (i.e., step)  $l = 0$  and
  observe state  $s_l$ 
  while  $s_l$  is not terminal and  $l < L$  do
    Execution Loop
      Exploration Loop  $k \in K$  do
        Compute  $A_i(t)$  by Eq: (1)
        Compute  $P_i(t)$  by Eq: (2)
        Compute  $V_k(t)$  by Eq: (4)
      end
      Update similarity by Eq: (7)
      using each data point in  $D$ 
      Predict student action  $a$  by
         $k_l \in \arg \max_{k \in K} V_k(t)$ 
    end
    Observe student action  $a$ , observe
     $s_{l+1}$ , and student feedback
    outcome  $u_{l+1}$ 
    Store  $t$  instance in  $\mathcal{M}$ 
  end
until task stopping condition

```

Algorithm 1: Pseudo Code of Instance-Based Learning Cosine Similarity Update

6 Case Study of IBIS: Individuals in Phishing Email Education Dataset

Previous comparisons of similarity metrics and human participant behavior compared the average of human performance. To highlight the benefits of the IBIS method, we replicate these calculations with one individual from the experiment. Here, the individual similarity of phishing and ham emails is based only on a single individuals categorization, confidence, and reaction time in their judgement. These graphs are shown for illustration with the average accuracy of logistic regression of similarity metrics predicting individual participant similarity metrics reported in table 1.

While the previous comparisons of embedding similarity metrics were all reasonably reflective of the average of all human participants across the entire dataset, they do not necessarily correspond to individual participants as closely. To demonstrate this, we plot 5 randomly selected participants’ individual metrics of similarity for the limited emails they observed in Figure 5. Here, the individual

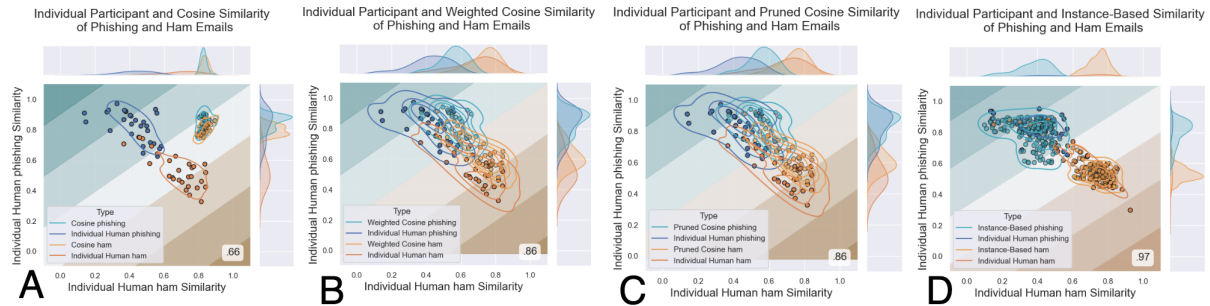


Figure 5: A comparison of the four similarity metrics under comparison using data from a single individual participant. These are shown for illustrative purposes with averages for each participant listed in Table 1 A: Cosine similarity compared to human subjective similarity. B: Weighted cosine similarity compared to human subjective similarity. C: Pruned cosine similarity compared to human subjective similarity. D: IBIS similarity compared to human subjective similarity.

441 similarity of phishing and ham emails is based only
 442 on a single individuals categorization, confidence,
 443 and reaction time in their judgement. From this we
 444 can see that there is a large discrepancy between
 445 the aggregated weighted cosine similarity and each
 446 of the four individual participants.

447 The accuracy of the logistic regression of the em-
 448 beddings for the vanilla cosine similarity for this
 449 example participant is 0.66. The same value for
 450 both the weighted cosine and pruned cosine method
 451 for this participant is 0.86. Meanwhile, the IBIS
 452 metric gives an accuracy of the logistic regression
 453 of 0.96. This is approaching the original accuracy
 454 of the two best performing cosine similarity met-
 455 rics (weighting and pruning) when using the entire
 456 dataset of human participant performance.

457 An important difference between these four
 458 methods is that only the IBIS method can com-
 459 pare emails that were not originally presented to
 460 an individual, meaning there are more embedding
 461 similarities used in the logistic regression. Note
 462 the smaller number of similarities performed in the
 463 Cosine, Weighted Cosine, and Pruned Cosine con-
 464 ditions in the first three columns of Figure 5, which
 465 is due to the limited number of emails shown to
 466 each individual participant. Meanwhile, the IBIS
 467 method has a larger sample of emails to draw from
 468 since it makes predictions of individual participant
 469 behavior on emails that they were never presented
 470 with.

471 This comparison demonstrates the clear benefits
 472 of using a cognitively inspired method of model-
 473 ing human participant decisions making that takes
 474 into account biases and cognitive constraints. The
 475 results is a prediction of behavior that can accu-
 476 rately fill in the gaps of unseen elements of the

477 dataset that have not been observed by a partici-
 478 pant. This method more accurately predicts the
 479 subjective similarity of participants as measured
 480 by categorization, confidence, and reaction time.
 481 Importantly, this is done while initially limiting
 482 the cognitive model to observing a single decision
 483 made by these participants, and increasing this data
 484 as the participant makes more decisions.

7 Predicting Human Categorization 485

486 To evaluate the usefulness of these previously listed
 487 metrics of semantic similarity, we employ another
 488 IBL model to make predictions of participant cate-
 489 gorizations of phishing emails as being either safe
 490 or dangerous. This is a highly complex task that
 491 involves making judgements about subjective qual-
 492 ities like suspiciousness, urgency, or plausibility,
 493 as well as objective qualities like whether the email
 494 sender matches one listed in the body or whether a
 495 link url is mismatched from its text.

496 In the above examples of using different methods
 497 to predict the subjective similarity of human partici-
 498 pant behavior, the entire dataset of decisions from
 499 one individual was used. However, when provid-
 500 ing educational feedback the number of data points
 501 from each participant begins at 0 and progresses
 502 through to the full amount of decisions collected
 503 from that participant. This is a significantly more
 504 challenging problem, as human participant deci-
 505 sions can be poor at the beginning of educational
 506 examples and potentially increase in quality dra-
 507 matically through educational feedback.

508 To compare the methods discussed in this pa-
 509 per, as well as our proposed Instance-Based cosine
 510 similarity weighting approach, we evaluate the ac-
 511 curacy of logistic regressions as they are formed

Method	Similarity to Average Humans	Similarity to Individual Humans	Human Behavior IBL Prediction Accuracy
Cosine Similarity (Park et al., 2020)	0.48	0.60±0.2	0.80±0.1
Embedding Weighting (Li and Han, 2013)	0.97	0.86±0.1	0.81±0.04
Embedding Pruning (Manrique et al., 2023)	0.97	0.86±0.04	0.82±0.08
IBIS (proposed)	0.97	0.93±0.04	0.87±0.05

Table 1: Comparison of the three previously described methods in their similarity to human behavior. Similarity to average humans is performed across the entire dataset of human judgements. Similarity to individuals and IBL prediction accuracy are both done for each individual participant. Reported values are means \pm standard deviations.

from individual participants behavior. This is done by comparing the regression accuracy in predicting the single next decision made by a human participant while fitting the measure of their subjective similarity from all previous decisions that they have made.

The results from this comparison of the predictive accuracy of a separate IBL model that relies on different metrics of similarity when predicting human performance are shown in Table 1, and indicate that the IBIS method of calculating individualized subjective similarity of participants produces the best similarity metric for an IBL model when predicting human behavior. It is important to note that the IBL model predicting human behavior and the model that are estimating similarity are not the same, as the similarity estimate model needs to rely on a separate similarity metric.

8 Discussion

Many applications of LLMs are interested in tailoring use cases to individuals, even when little information is known about that individual. While many approaches of individualization have demonstrated success in producing outputs or representing information in an individualized manner, these have typically relied on advanced machine learning techniques. The method proposed in this work is relatively simple from a mathematical perspective, though there is a strength in its reliance on theories of cognition that underlie human learning and decision making. The result is a simple to understand and easy to implement method of calculating similarities of unseen documents using a cognitive model, which can augment datasets that contain only a small number of decisions from a single user.

The specific application we investigated is somewhat unique in that it is based on training human

participants to make categorization judgements of textual information of one of two categories. However, we believe that the general method described, of augmenting subjective similarity metrics with predicted decisions from a cognitive model, could be applied onto various other scenarios.

For instance, in visual learning settings Variational Autoencoders have been integrated with cognitive models to predict human utility learning of abstract visual information (Malloy and Sims, 2024). This task involved visual queues with associated utilities taken from a large dataset of hundreds of possible abstract visual images in the form of jars of differently colored marbles. The same method of determining subjective similarity could be applied onto this visual utility learning task.

Overall, the results in this work demonstrate the usefulness of cognitive models in serving as digital twins to human participants. Leveraging these models and integrating their results into Large Language Model techniques has the potential to make measurements from these models more cognitively grounded. While there are existing methods of incorporating human behavior through the use of large datasets collected from many participants, these do not necessarily account for individual biases and constraints. The method proposed in this work takes these features of human learning and decision making into account in developing individualized metrics of similarity.

Limitations

The task presented in this work of predicting whether an email is phishing or ham relies heavily on a small number of features within the email. Namely, if an email contains a link that redirects to a nefarious website, or requests personal information, then it should be labelled as phishing. While students rely on many queues to make their judge-

ments, the true categorization task is in reality simple. Future work in the area of learning subjective similarity metrics should expand into domains with more categories, and more complex and abstract categories.

Ethics Statement

The model proposed in this work, as well as the dataset introduced, involves an educational setting and thus introduces significant ethical concerns. One of the main concerns of the use of LLMs in educational settings is the potential for biases present in LLMs that negatively impact students of a specific ethnic, cultural, or racial background. This potential concern is mitigated in this work because of the specific educational setting, in detecting phishing emails, which are designed by the original cybersecurity experts to be applicable to a wide range of end users. However, the application of this approach outside of the setting used in this work should take care in ensuring that the method of calculating the similarity of educational examples shown to students not be biased. While this is an inherent concern in the use of LLMs in education, our proposed approach of using more individualized metrics of similarity can hopefully reduce the likelihood of LLM biases negatively impacting student education. This is because our proposed model is based on individual past experiences and biases when calculating subjective similarity.

References

Palvi Aggarwal, Omkar Thakoor, Shahin Jabbari, Edward A Cranford, Christian Lebiere, Milind Tambe, and Cleotilde Gonzalez. 2022. Designing effective masking strategies for cyberdefense through human experimentation and cognitive models. *Computers & Security*, 117:102671.

Shivani G Aithal, Abishek B Rao, and Sanjay Singh. 2021. Automatic question-answer pairs generation and question similarity mechanism in question answering system. *Applied Intelligence*, pages 1–14.

Asim Ansari, Skander Essegaier, and Rajeev Kohli. 2000. Internet recommendation systems.

Luís Borges, Bruno Martins, and Pável Calado. 2019. Combining similarity features and deep representation learning for stance detection in the context of checking fake news. *Journal of Data and Information Quality (JDIQ)*, 11(3):1–26.

Edward A Cranford, Cleotilde Gonzalez, Palvi Aggarwal, Sarah Cooney, Milind Tambe, and Christian

Lebiere. 2020. Toward personalized deceptive signaling for cyber defense using cognitive models. *Topics in Cognitive Science*, 12(3):992–1011.

Edward A Cranford, Christian Lebiere, Prashanth Rajivan, Palvi Aggarwal, and Cleotilde Gonzalez. 2019. Modeling cognitive dynamics in end-user response to phishing emails. *Proceedings of the 17th ICCM*.

Cleotilde Gonzalez and Varun Dutt. 2011. Instance-based learning: integrating sampling and repeated decisions from experience. *Psychological review*, 118(4):523.

Cleotilde Gonzalez, Javier F Lerch, and Christian Lebiere. 2003. Instance-based learning in dynamic decision making. *Cognitive Science*, 27(4):591–635.

Ahmed Yar Khan, Rabia Latif, Seemab Latif, Shahzaib Tahir, Gohar Batool, and Tanzila Saba. 2019. Malicious insider attack detection in iots using data analytics. *IEEE Access*, 8:11743–11753.

Harsh Khatter, Nishtha Goel, Naina Gupta, and Muskan Gulati. 2021. Movie recommendation system using cosine similarity with sentiment analysis. In *2021 Third International Conference on Inventive Research in Computing Applications (ICIRCA)*, pages 597–603. IEEE.

Alfirna Rizqi Lahitani, Adhistya Erna Permanasari, and Noor Akhmad Setiawan. 2016. Cosine similarity to determine similarity measure: Study case in online essay assessment. In *2016 4th International conference on cyber and IT service management*, pages 1–6. IEEE.

Tomás Lejarraga, Varun Dutt, and Cleotilde Gonzalez. 2012. Instance-based learning: A general model of repeated binary choice. *Journal of Behavioral Decision Making*, 25(2):143–153.

Baoli Li and Liping Han. 2013. Distance weighted cosine similarity measure for text classification. In *Intelligent Data Engineering and Automated Learning—IDEAL 2013: 14th International Conference, IDEAL 2013, Hefei, China, October 20–23, 2013. Proceedings 14*, pages 611–618. Springer.

Tyler Malloy and Cleotilde Gonzalez. 2024. Applying generative artificial intelligence to cognitive models of decision making. *Frontiers in Psychology*, 15:1387948.

Tyler Malloy and Chris R Sims. 2024. Efficient visual representations for learning and decision making. *Psychological review*, in press.

Natalia Flechas Manrique, Wanqian Bao, Aurelie Herbelot, and Uri Hasson. 2023. Enhancing interpretability using human similarity judgements to prune word embeddings. *arXiv preprint arXiv:2310.10262*.

Routa Moussaileb, Nora Cuppens, Jean-Louis Lanet, and Hélène Le Boudier. 2021. A survey on windows-based ransomware taxonomy and detection mechanisms. *ACM Computing Surveys (CSUR)*, 54(6):1–36.

693 Shaimaa M Nafea, François Siewe, and Ying He. 2019.
694 A novel algorithm for course learning object rec-
695 ommendation based on student learning styles. In
696 *2019 International Conference on Innovative Trends
697 in Computer Engineering (ITCE)*, pages 192–201.
698 IEEE.

699 Thuy Ngoc Nguyen and Cleotilde Gonzalez. 2022. The-
700 ory of mind from observation in cognitive models
701 and humans. *Topics in Cognitive Science*, 14(4):665–
702 686.

703 Kwangil Park, June Seok Hong, and Wooju Kim. 2020.
704 A methodology combining cosine similarity with
705 classifier for text classification. *Applied Artificial
706 Intelligence*, 34(5):396–411.

707 Rajvardhan Patil, Sorio Boit, Venkat Gudivada, and
708 Jagadeesh Nandigam. 2023. A survey of text repre-
709 sentation and embedding techniques in nlp. *IEEE
710 Access*.

711 Mansoore Shojaei and Hassan Saneifar. 2021. Mfsr: A
712 novel multi-level fuzzy similarity measure for recom-
713 mender systems. *Expert Systems with Applications*,
714 177:114969.

715 Xuansheng Wu, Xinyu He, Tianming Liu, Ninghao Liu,
716 and Xiaoming Zhai. 2023. Matching exemplar as
717 next sentence prediction (mensp): Zero-shot prompt
718 learning for automatic scoring in science education.
719 In *International conference on artificial intelligence
720 in education*, pages 401–413. Springer.

721 Peipei Xia, Li Zhang, and Fanzhang Li. 2015. Learning
722 similarity with cosine similarity ensemble. *Informa-
723 tion sciences*, 307:39–52.

724 Kaitlyn Zhou, Kawin Ethayarajh, Dallas Card, and Dan
725 Jurafsky. 2022. Problems with cosine as a measure
726 of embedding similarity for high frequency words.
727 *arXiv preprint arXiv:2205.05092*.