

ACT-R Workshop Individual Differences 2: Personalized Model-driven Interventions for Decisions from Experience

Presented by Drew Cranford

Personalizing human-machine interactions

- Human-machine interactions are typically static
 - System policies and algorithms often tailored to populations, or average human behavior
 - Sometimes based on erroneous assumptions/models of human behavior
 - E.g., that humans make perfectly rational decisions.
- Can cognitive architectures be used to personalize human-machine interactions?
 - Yes
 - Prime examples from intelligent tutors¹
 - Can be cumbersome to model complete scope of task



¹http://act-r.psy.cmu.edu/category/other/intelligenttutoring-systems/

Decisions from Experience

- Within a task, there are many instances of decisions from experience.
 - We can leverage the very powerful modeling methodology of Instance-Based Learning to accurately model decisions from experience
- Instance-Based Learning Theory¹
 - Decisions made by generalizing across similar past experiences
 - Contextual features
 - Action
 - Outcome/Reward (utility)
 - ACT-R Blending mechanism



Insider Attack Game (IAG)

MURI: Cyber Deception



Cranford, E. A., Gonzalez, C., Aggarwal, P., Cooney, S., Tambe, M., & Lebiere, C. (2019). Towards personalized deceptive signaling for cyber defense using cognitive models. *ICCM 2019*. Montreal, CA.

4

IBL ACT-R Model Logic



Aligning Individual Human and Model Runs

- Humans behave differently from one another
 - They learn, and are adaptive
 - Have unique cognitive biases that arise from unique experiences
 - (e.g., confirmation bias)





Aligning Individual Human and Model Runs

- Model-tracing/Knowledge-tracing¹
 - Run model alongside human
 - Predict human decisions
 →feed data to system
 - Adapt system
 →then human makes decision
 - Add/modify chunks to match human decisions



Aligning Individual Human and Model Runs

• 2 chunks changed:

Probability of Attack (Model)

1.00-

0.75-

0.50-

0.25-

0.00

0.00

- Ground truth decision
 - Observed behavior

 $r^2 = 0.95$

0.25

0.50

Probability of Attack (Human)

0.75

1.00

- Expectation (e.g., confirmation bias)
 - e.g., if model expects negative outcome but human attacked, then *infer* that the expected outcome was the reward value



Cranford, E. A., Gonzalez, C., Aggarwal, P., Cooney, S., Tambe, M., & Lebiere, C. (2019). Towards personalized deceptive signaling for cyber defense using cognitive models. *ICCM 2019*. Montreal, CA.

Cognitive Signaling – Costs and Benefits

- Humans lose trust in the signal when catching it being deceptive
- Restore trust with blocks of truthful signals; but cost of giving free pass
- Optimize tradeoff between short-term cost and long-term gains:

• **Current probability of attack:**
$$P_{est}^{now}(A|S) = \frac{\sum_{i}^{wins} t_i^{-d} + \sum_{j}^{losses} t_j^{-d}}{\sum_{i}^{wins} t_i^{-d} + \sum_{j}^{losses} t_j^{-d} + \sum_{k}^{draws} t_k^{-d}}$$

Assumed probability of attack given signal:

0

$$P_{ass}^{now}(A|S) = \frac{\sum_{i}^{wins} t_i^{-d}}{\sum_{i}^{wins} t_i^{-d} + \sum_{j}^{losses} t_j^{-d}}$$

• Expected number additional losses: $1/3 * b * P_{est}^{now}(A|S)$

Issue truthful block if costs < benefits:
$$\frac{1}{3} * b * [1 - P_{est}^{now}(A|S)] < \alpha * r * [P_{ass}^{now}(A|S) - P_{ass}^{then}(A|S)]$$

Cranford, E. A., Gonzalez, C., Aggarwal, P., Cooney, S., Tambe, M., & Lebiere, C. (2020). Adaptive cyber deception: Cognitively-informed signaling for cyber defense. *HICSS 53*, January 2020 (pp. 1885-1894). Maui, HI.

Model predictions vs human behavior



Aggressive attackers report ignoring the signal when making the attack decision



Modifying representation







What's next?

- Future work aimed at optimizing adaptive signaling scheme
 - Lessons from reinforcement learning
 - Agent = signal
 - Actions = present; withhold
 - Goal = minimize expected outcome generated via blending
- For some participants, no matter how the signaling scheme changes, they ignore the signal and therefore behavior will never change.
 - Adapt coverage as well?
 - Based on target selection preferences
- To accurately model some participants, we need to adapt the model beyond adjusting declarative memory.
 - Adapt model representation of chunks?

What other kinds of info can we use to adapt models?

- Other patterns of behavior (e.g., selection preferences)
- Underlying cognitive states (proxy for mental model?)
 - Distribution of instances in memory
 - Activation of chunks
 - Probability of retrieving instances
 - Feature salience
- If we can determine a feature is not important, we can exclude that feature from the model representation.
 - More accurate predictions
 - Cognitive Salience
 - Degree of influence of individual features

Cognitive Salience

• Take derivative of blending equation with respect to each feature³



³Somers, S., Mitsopoulos, K., Lebiere, C., & Thomson, R. (2019). Cognitive-Level Salience for Explainable Artificial Intelligence. In *Proceedings of the 17th Annual Meeting of the International Conference on Cognitive Modeling*. Montreal, CA.

Personalized coverage using salience



- This information could be used to adapt coverage
- Also, can be used to adapt model to make more accurate predictions

https://virtual.mathpsych.org/presentation/195

Personalized signaling using salience



- This information could be used to adapt signaling
- Also, can be used to adapt model to make more accurate predictions

Conclusions

- Combination of model/knowledge-tracing techniques and IBL models provides efficient method to personalize human-machine interactions
 - Very accurate predictions of individuals
 - Very little data (experience) needed to make them
 - Efficient application
 - General decision process; few production rules
 - Typically only need to modify declarative memory
 - Architecture provides invaluable information to personalize systems
 - Explanations of behavior in addition to decision predictions
- Future research aimed at addressing scalability issues
 - ACT-UP/pyACT-UP
 - David Reitter, Christian Lebiere¹; Don Morrison²
 - Vectorized memory for ACT-R
 - Matthew Kelly³
 - Vectorized IBL
 - Konstantinos Mitsopoulos

¹Reitter, D., & Lebiere, C. (2010). Accountable Modeling in ACT-UP, a Scalable, Rapid-Prototyping ACT-R Implementation. In Proceedings of the 2010 International Conference on Cognitive Modeling. Philadelphia, PA.
²https://bitbucket.org/dfmorrison/pyactup

³Kelly, M. A., Arora, N., West R. L., & Reitter, D. *High-dimensional vector spaces as the architecture of cognition*. Manuscript in submission as of September 2019. doi: http://dx.doi.org/10.31234/osf.io/ryvg2

PIs & Collaborators

Christian Lebiere - PI Coty Gonzalez - PI Milind Tambe - Pl Palvi Aggarwal Sarah Cooney **Sterling Somers** Konstantinos Mitsopolous Don Morrison

Questions?





SENTED HERE THE RESEARCH PRE WAS LARG D BY THE ARMY RESEA FICE AND ACCOMPI GRANT NUMBER W911NF-17-1-0370.

August 31, 2020

MathPsych/ICCM 2020

ACT-R Workshop