ACT-R PGSS 2016

WRIGHT STATE

Learning to trust and trusting to learn

Ion Juvina, Othalia Larue, Michael G. Collins, & Peter Crowe

Department of Psychology Wright State University Dayton, OH

Outline

- Broader architectural issues
- A published model of learned trust
 - Validation studies
 - Critique: strengths and limitations
- A revised model

*

WRIGHT STATE

- Validation
 - External validity
 - Predictions
- Conclusion and future work



WRIGHT STATE UNIVERSITY Broader architectural issues

- ACT-R is best positioned to inform the rational and social bands
 - Develop process models of decision making and strategic interaction
 - Lasting, repeated, cooperative or adversarial
 - Involving rational agents balancing multiple motives, constraints etc.
 - Facilitate incremental theory building and integration
 - Use cognitive models as intelligent agents/robots
 - e.g., companion robots



WRIGHT STATE Broader architectural issues

- · Representing agents vs. representing inanimate entities
 - Theory-of-mind
 - Impression formation
 - Trust
 - Human-machine interaction
 - Anthropomorphic features

WRIGHT STATE **Objective & approach**

- Develop a unified theory of learned trust ٠ - grounded in general principles of human cognition
 - A priori predictions

*

- Predictions generated by a computational model <u>before</u> a human study is run (ex ante)
- Model developed based on theory, literature, or prior studies - Generate predictions for
- New task, new conditions / manipulations
- Assess validity of model under new conditions
- Study design and setup are identical for model simulations and human data collection
- Simulation data are as rich and fine-grained as the human data
- Revise model
- Model -> predictions -> experiments -> model ...

WRIGHT STATE UNIVERSITY

The published model

Juvina, I., Saleem, M., Martin, J.M., Gonzalez, C., & Lebiere, C. (2013). Reciprocal trust mediates deep transfer of learning between games of strategic interaction. Organizational Behavior and Human Decision Processes. 120(2): 206-215.

Juvina, I., Lebiere, C., & Gonzalez, C. (2015). Modeling trust dynamics in strategic interaction. Journal of applied research in memory and cognition. 4(3): 197-211.







	WRIGHT STATE UNIVERSITY				
Determines which reward function to use					
	Trust Tr				
	inv				
	+ ·				
	+ -				
12					

WRIGHT STATE UNIVERSITY MOdel comparison				
Model	Training	z sample	Testing	sample
	Correlation	RMSD	Correlation	RMSD
Trust	0.87	0.09	0.81	0.11
P1 _n	0.66	0.19	0.65	0.19
$P1_n - P2_n$	0.01	0.46	- 0.001	0.46
$P1_n + P2_n$	0.65	0.29	0.61	0.30
$P1_n + P2_n - P2_{n-1}$	0.66	0.25	0.59	0.27
P2 _n	0.62	0.34	0.62	0.34
				13

WRIGHT STATE UNIVERSITY MOdel validation
 Validation study 1 Same games (PD & CG) A wider range of conditions
A priori model predictions (published)
Large set of numan data (320 participants)
Collins, M.G., Juvina, I., & Gluck, K.A. (2016). Cognitive Model of Trust Dynamics Predicts Human Behavior within and between Two Games of Strategic Interaction with Computerized Confederate Agents. <i>Frontiers in Psychology, Section Cognitive</i> <i>Science</i> 7:49.
14

WRIGHT STATE UNIVERSITY	WRIGHT STATE Setup: what was different
 Trait trust = trust propensity, dispositional trust State trust = specific to a (repeated) interaction Trust necessity = need to invest in trust development 	 Number of rounds 50 instead of 200 Payoff matrix (-4,-1,1,4) instead of (-10,-1,1,10) Human vs. confederate agent instead of human vs. human Participants Amazon's Mechanical Turk workers Instead of CMU undergraduates Order conditions (PDPD, PDCG, CGPD, CGCG) instead of (PDCG, CGPD) Strategy of confederate agent Tit-for-tat

_

WRIGHT STATE UNIVERSITY Study design					
	Gan	ne Order Confederate Agent			
	First Game	Second Game	Strategy Trustworthiness		
	PD	PD	T4T	High Trustworthiness	
	PD	CG	T4T	High Trustworthiness	
	CG	PD	T4T	High Trustworthiness	
	CG	CG	T4T	High Trustworthiness	
	PD	PD	T4T	Low Trustworthiness	
	PD	CG	T4T	Low Trustworthiness	
	CG	PD	T4T	Low Trustworthiness	
	CG	CG	T4T	Low Trustworthiness	
	PD	PD	PT4T	High Trustworthiness	
	PD	CG	PT4T	High Trustworthiness	
	CG	PD	PT4T	High Trustworthiness	
	CG	CG	PT4T	High Trustworthiness	
	PD	PD	PT4T	Low Trustworthiness	
	PD	CG	PT4T	Low Trustworthiness	
	CG	PD	PT4T	Low Trustworthiness	
	CG	CG	PT4T	Low Trustworthiness	

WRIGHT STATE Trait and state trust scales

• Trait trust scale (24 items):

- I generally trust other people unless they give me a reason not to.
- One is better off being cautious when dealing with strangers until they have provided evidence that they are trustworthy.
- State trust scale (14 items):
 - I believe that the other player wants to help me to make a good amount of payoff in this game.
 - The other player can be trusted.
 - I would not let the other player have any influence over my payoff.





Study 1 Conclusions

- · Majority of model predictions was observed
- Deviations from predictions
 - Humans cooperate

.

WRIGHT STATE

- More against low-trustworthy agents than predicted
- Less against high-trustworthy agents than predicted
- Different setups, different populations
- Relationship between trait trust and state trust – Both trait trust and state trust are learned.
- Trust asymmetry

WRIGHT STATE Model validation cont'd.

• Validation study 2

21

- Same design as validation study 1 (16 conditions)
- Added counterpart-change manipulation (32 conditions)
- A priori model predictions (published)
- Larger set of human data (640 participants) – Not published yet.





WRIGHT STATE UNIVERSITY

- Strengths (skip for now)
- Limitations

WRIGHT STATE UNIVERSITY					
No link between trait and state trust					
– State trust starts at zero?					
 Initial trust is high (McKnight, Cummi Chervany, 1998) 	ngs, &				
 Trait trust influences state trust (Berg McCabe, 1995; Dirks & Ferrin, 2001) 	ı, Dickhaut, &				
- State trust influences trait trust (Collins	et al., 2016)				
 Novel finding 					
 Relevant for interactions with multiple sequence 	e trustees in				
	26				

WRIGHT STATE Limitations (cont'd)

- Learning equation was linear
 - $ST_t = ST_{t-1} + PET_t$
- Did not explain
 - Trust asymmetry (Slovic, 1993)
 - Early evidence more important than late evidence (Lount et al., 2008)
- Learning equations tend to be power functions (Newell & Rosenbloom, 1981)
- · No link between cognitive ability and learned trust
 - Cognitive ability is one of the best predictors of learning and performance

Goal: Overcome limitations of the published model

- Expand the model's scope of applicability
- External validity

Juvina, I., Collins, M.G., Larue, O., & de Melo, C. (2016). *Toward a unified theory of learned trust*. Paper presented at the International Conference on Cognitive Modeling, State College, PA.

WRIGHT STATE UNIVERSITY

- Initial state trust = trait trust
 - Modulo some incidental learning (emotions, gaze etc.)
- Trait trust = f(prior state trusts)
 (Collins et al., 2016)
- Trait trust deviation = TT_end TT_start
- Cognitive ability = accuracy of judgments of trustworthiness and trust necessity

WRITE TARE Evised trust learning equation

$$ST_t = ST_{t-1}^a + PET_t - b * TTD$$

- ST = state trust
- a = power exponent; a < 1
- PET = perceived evidence of trustworthiness
- TTD = trust propensity deviation
- b = perception bias







- Test that revised model:
 - Maintains published model's strengths
 - Explains other (seemingly unrelated) findings
 - Makes novel predictions





Jutcome	Publish	ed model	Revis	ed model
-	Trust	Invest	Trust	Invest
CC	3	NA	6	NA
CD	-10	-1	-7	-1
DC	10	NA	9	NA
DD	-1	.18	-1*	.18*
ublished evised m	model: r	= .89, RMS .90. RMSD	SD = .09) = .07	











WRIGHT STATE UNIVERSITY	Fit stats Lount			
• Revised model: <i>r</i> = .99, <i>RMSD</i> = .33				
		42		





Fit stats De Melo

• Revised model: *r* = .86, *RMSD* = .11

*

WRIGHT STATE

Conclusions Conclusions Model has potential to unify theories on learned trust – Cumulates learning from History of prior interactions (trait trust) Evidence of trustworthiness and trust necessity in current interaction (Incidental learning from facial expressions)

Predictions Predictions Trust decay Evidence? Cognitive ability influences state and trait trust Lyons, Stokes, & Schneider, 2011; Sturgis, Read, & Allum, 2010;

- Yamagishi, Kikuchi, & Kosugi, 1999



WRIGHT STAWORK IN progress and future work

Empirical

- Using models as autonomous agents
- Study in progress
- Study to test the trust decay prediction
- Modeling:
 - How state trust influences trait trust
 - TT = f(ST)
 - Dynamics of trait trust and perception bias
 - b = f(TTD)

Collaborators & funding

- Kevin Gluck
- Celso de Melo
- Randall Green

The work presented here was supported by The Air Force Office of Scientific Research grant number FA9550-14-1-0206 to lon Juvina and the Oak Ridge Institute for Science and Education (ORISE) who supported this research by appointing Michael Collins to the Student Research Participant Program at the U.S. Air Force Research Laboratory (USAFRL), 711th Human Performance Wing, Human Effectiveness Directorate, Warfighter Readiness Research Division, Cognitive Models and Agents Branch administered by the ORISE through an interagency agreement between the U.S. Department of Energy and USAFRL.

WRIGHT STATE UNIVERSITY	Questions?	