

Political scientists identify two categories of problems that impede trust: *information problems* and *commitment problems* [5]. They are traditionally applied to the prisoner’s dilemma, but they play a critical role in understanding problems of bargaining in war [13], iterated games [1], and economic policy [16].

For instance, the current U.S.–China trade war presents a setting in which two parties converge on an ostensibly suboptimal policy (economic sanctions) because of these problems. First, each party lacks complete information on the long-term objectives of the other (e.g. geopolitical / economic agendas). Second, each party cannot be sure that the other will commit to policy changes (e.g. reduction in trade deficit, or improvement in Chinese market practices). Modelling the behavior of agents with information theory might play a role in building AI we can trust.

There are three core facets of “information-theoretic trust”: *compression*, *communication*, and *inference*. While these are standard topics in introductory texts [3][14], it is worth expanding on the subtle connections between information theory and trust. I will elaborate on these three topics in the next few paragraphs.

First, compression plays a critical role in safe interaction with an ensemble of agents. Codes have played an essential role in history, e.g. see the prominent role of the Enigma machine in World War II [2]. Secure compression requires understanding how much we can safely encode in a message.

Shannon’s source coding theorem states that N i.i.d. random variables with entropy $H(X)$ can be compressed into more than $NH(X)$ bits with negligible risk of information loss as $N \rightarrow \infty$ [17][14]. Conversely, if they are compressed into fewer than $NH(X)$ bits, it is virtually certain that information will be lost. The language of compression is a powerful framework to analyze dimensionality reduction models like principal component analysis [7] and autoencoders [4].

Second, safely interacting with AI requires trust in our communication channels. While the source coding theorem provides a useful framework to analyze the lossless case, in practice, we often deal with noisy channels. “Nuclear close calls” present useful case studies to examine in history, wherein nuclear weapons were “nearly” deployed by a legitimate authority [18]. In these settings, state actors are unsure of the intentions of others. The Norwegian rocket launch of 1995, described as a “poster child for nuclear dangers” (Tetra [18]), involved Norwegian and American scientists launching a rocket to study weather data. While information about this launch was sent to Moscow, before it reached the authorities, the launch was interpreted as a sign of a possible adversarial strike.

Shannon’s noisy-channel coding theorem formalizes the notion of channel capacity, a measure of how much we can safely transmit over a noisy channel. We first review the *mutual information*, which provides a measure of how much information is communicated over a channel. If X, Y are random variables representing the input and output to a noisy channel, then the mutual information quantifies how much information the output conveys about the input: $I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X)$. The *channel capacity* of a channel Q can be computed as $C(Q) = \max_{P_X} I(X; Y)$, where the maximization is over all possible input ensembles P_X . Intuitively, $C(Q)$ represents the number of bits that can be sent over a channel with arbitrarily low decoding error for the recipient.

Finally, trusting AI requires us to model and infer the behavior of complex agents. We might initially ask which priors are most accurate in settings of incomplete information. There are many potential answers to this question, but information theory suggests that the “maximum-entropy distribution” makes the fewest assumptions [6] [14].

One can derive the familiar probability distributions using this idea. With fixed variance, the Gaussian is the maximum entropy prior; with fixed mean on $(0, \infty)$, the exponential distribution is the maximum entropy prior. There are an entire class of methods, “maximum entropy methods,” with applications in time series forecasting, combinatorics, and statistical analysis more broadly

[9][10].

While model interpretability techniques [8][12][15] have shed light on how to understand agents, information-theoretic models are still powerful tools we can apply to model the multi-agent case. Build a future with safe, trustworthy AI may require deep inquiry and understanding of the information-theoretic toolbox.

Acknowledgments

Thanks to Yuval Wigderson and Michael Swerdlow for reading drafts of this essay.

References

- [1] Robert Axelrod and William D Hamilton. The evolution of cooperation. *science*, 211(4489):1390–1396, 1981.
- [2] Bernhelm Booß-Bavnbek and Jens Høyrup. *Mathematics and war*. Springer, 2003.
- [3] Thomas M Cover and Joy A Thomas. *Elements of information theory*. John Wiley & Sons, 2012.
- [4] Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.
- [5] Jeffrey A Frieden, David A Lake, and Kenneth A Schultz. World politics. interests, interactions. *Institutions*. WW Norton & Company, New York, 2010.
- [6] Adithya C Ganesh. The principle of maximum entropy. *Stanford mathematics department: directed reading program*, 2019.
- [7] Bernhard C Geiger and Gernot Kubin. Relative information loss in the pca. In *2012 IEEE Information Theory Workshop*, pages 562–566. IEEE, 2012.
- [8] Fred Matthew Hohman, Minsuk Kahng, Robert Pienta, and Duen Horng Chau. Visual analytics in deep learning: An interrogative survey for the next frontiers. *IEEE transactions on visualization and computer graphics*, 2018.
- [9] Edwin T Jaynes. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(9):939–952, 1982.
- [10] Edwin T Jaynes. *Probability theory: The logic of science*. Cambridge university press, 2003.
- [11] Hyunjik Kim and Andriy Mnih. Disentangling by factorising. *arXiv preprint arXiv:1802.05983*, 2018.
- [12] Pang Wei Koh and Percy Liang. Understanding black-box predictions via influence functions. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1885–1894. JMLR. org, 2017.
- [13] David A Lake. Two cheers for bargaining theory: Assessing rationalist explanations of the iraq war. *International Security*, 35(3):7–52, 2010.
- [14] David JC MacKay. *Information theory, inference and learning algorithms*. Cambridge university press, 2003.
- [15] Chris Olah, Alexander Mordvintsev, and Ludwig Schubert. Feature visualization. *Distill*, 2(11):e7, 2017.
- [16] Arvind Panagariya. International trade. *Foreign Policy*, pages 20–28, 2003.
- [17] Claude Elwood Shannon. A mathematical theory of communication. *Bell system technical journal*, 27(3):379–423, 1948.

- [18] Bruno Tertrais. “on the brink”—really? revisiting nuclear close calls since 1945. *The Washington Quarterly*, 40(2):51–66, 2017.