# STATEMENT OF PURPOSE
Hadi Khalaf

I am interested in using tools from information theory and statistics to study problems in distributed optimization and learning. With larger real-world datasets and systems, we need secure machine learning models with theoretical guarantees on their performance. Throughout my undergraduate studies, I had the opportunity to think about such models from the standpoint of three questions:

1. How can we produce models "complex" enough that they express a behavior of interest?

2. How can we produce estimates that work well in the presence of outliers?

3. How can we produce estimates under privacy constraints in distributed data settings?

My research at Harvard University and the American University of Beirut has focused on learning from complex systems. This work was also informed by my technical experience designing ML-powered chatbots for community-based enterprises in Lebanon. The challenges that arise from weak infrastructure call for human-centered solutions and usable frameworks. I hope my research can help expand existing efforts to create such open-source tools and, in turn, democratize access to information.

## When do models become complex?

Information theory has led to valuable tools for learning from data. In the information bottleneck (IB) problem, we want a compressed version of our observation $X$ that retains as much information as possible about a given target $Y$ [1]. While an exact solution can be found in the case of Gaussian distributions [2], there is little work on IB for general distributions to the best of my knowledge. This is critical for models that use IB as a form of regularization.

Hence, I reached out to Prof. Ibrahim Abou-Faycal, and we derived an analytical solution under strict regularity conditions. One issue we faced was that our solution, under some conditions, would compress $X$ to a single codeword. First, I thought it was surprising. However, I noticed a phase transition consistent with existing literature: the optimal model gradually prioritizes informativeness over compression. This is one example of the tradeoff between compression and expressivity. I am interested in understanding when a model transitions to an informative structure and extending this characterization of complexity to other systems. This may allow us to understand how compression can help a model generalize.

## What are 'good' estimates given an adversary?

Informative models allow us to learn from data and estimate unknown quantities. There is a growing interest in when estimation frameworks "break down" under an adversarial choice of data. This is especially important in higher dimensions where it is easier to generate adversarial data. I worked on partial identification with Prof. Elie Tamer at the Economics Department at Harvard University. It is an estimation framework that can be applied when

we cannot estimate model parameters using standard techniques. We focused on games where agents each choose to enter a market with a given probability.

We wanted to see how we can learn from existing data when market conditions or an agent's utility changes. To the best of our knowledge, there aren't tools that estimate bounds on entry probabilities in counterfactual cases. Hence, we developed novel algorithms for counterfactual learning using game-theoretic strategies like stochastic fictitious play and no-regrets. Instead of producing exact estimates, we found intervals where the true value might exist. One challenge was to evaluate confidence in these sets. I suggested fitting these intervals in an alpha shape, a generalization of convex hulls, to see where the intervals overlap most often. Intuitively, these should contain the true value with high probability. An interesting application is finding these bounds in the presence of outliers or censored data.

### How can agents learn together securely?

Next, we can look at the estimation paradigm across different agents. With the rapid increase in connected devices, fairness and privacy are vital to leverage collective intelligence and develop trustworthy and accessible models. For my final year project, I am designing the backend of a crowdsourcing platform for rental prices in Lebanon's capital to counter misinformation in market prices. This human-centered design requires us to study in practice how to limit the amount of shared data while maintaining a high-performance, reliable, and secure platform.

This project motivated my interest in federated learning as a scalable learning framework for decentralized networks. However, its performance degrades with data heterogeneity. I started working with Prof. Maher Nouiehed to find a sharpness-aware algorithm that aligns the global model with that of the clients. This is important since "flat" minimizers of the loss functions are suggested to promote the model's generalization ability [3]. To limit communication rounds and client drift, we need to "correct" how the parameters are updated by injecting some information about the global behavior. Based on experimental observations, I proposed adding a correction term that reflects how weights change. The next step is to see how the model behaves under differential privacy mechanisms.

## Future Plans

In the next few months, I will continue my research and complete more coursework in statistics. I will also lead a team to develop computational tools for rare-event simulation. Another central part would be teaching. I have been the head teaching assistant in the three largest courses the Engineering faculty offers. My interactions with the students and the teaching team have allowed me to experiment with different ways of communicating science. Hence, I plan to become a professor to advance my research and pass on the great education I have received.

Harvard would be a great place for me to pursue my Ph.D. I am interested in the interdisciplinary work of the Machine Learning Foundations Group. I want to explore Prof. Flavio

REFERENCES

Hadi Khalaf

Calmon's work on estimating information measures and studying utility-communication-privacy tradeoffs with an emphasis on differential privacy guarantees. With my background in statistics, I am interested in Prof. Yue Lu's theoretically grounded work on phase transitions that occur in high-dimensional inference. Similarly, I hope to help Prof. Cengiz Pehlevan build on his recent work on the grokking phenomenon and study learning dynamics from the viewpoint of the geometry of the loss landscape.

With my affinity for theoretical work and human-centered applications, I hope to contribute to the diverse community at Harvard as I explore new problems and move forward as a researcher.

# References

[1] N. Tishby, F. C. Pereira, and W. Bialek, "The information bottleneck method," 2000.

[2] G. Chechik, A. Globerson, N. Tishby, and Y. Weiss, "Information bottleneck for gaussian variables," *Journal of Machine Learning Research*, vol. 6, no. 6, pp. 165–188, 2005.

[3] P. Foret, A. Kleiner, H. Mobahi, and B. Neyshabur, "Sharpness-aware minimization for efficiently improving generalization," in *International Conference on Learning Representations*, 2021.

# Personal Statement

Hadi Khalaf

Beirut has been my home for the last two decades, and it has forced me to imagine a vibrant reality for myself and my surroundings. Despite political turmoil and economic freefall, there is always a commitment to move forward. This commitment usually starts in *counterspaces*: spaces that enable individuals to engage in collective. In my first year at the American University of Beirut, I joined the AUB Secular Club, the largest politically independent student club in the Middle East. What I love about the club is its potential: it was our only space for growth, the clash of ideas, and inclusivity. I became its president two years later.

Amidst a national crisis in 2021, the university's health insurance carrier limited coverage of psychiatry. I drafted a mental health reform proposal from the testimonies of more than two thousand students. It was the most comprehensive proposal to be submitted to the student council. It called for restructuring three offices at the university, extendliing access to mental health services, and training academic advisors on mental health aid. I led a year and a half of negotiations, and the university eventually reversed its decision. Under my presidency, the club also actively participated in the parliamentary elections. We encouraged students to vote, held discussions on campus and in local neighborhoods, and drafted policies to provide more financial support for Lebanon's only public university and equitable access to education for Lebanese students and refugees alike.

Achieving any new reality requires us to leverage collective efforts. I learned how to draw power from the everyday citizen by mobilizing people in my community. This motivated my interest in distributed models and human-centered designs: empathetic, innovative, and socially responsible systems that question what we assume to be true. I hope to have an active role in democratizing access to machine learning models because I have seen firsthand the impact of this technology in developing communities. In this direction, I am leading an effort to design chatbots that give market insight to small-business owners. Building on open-source frameworks, I designed these bots to be intuitive and user-friendly, lowering the barriers to resources typically reserved for larger corporations. I started with helping my father's house supplies store, and now I support nine businesses across Lebanon.

Graduate school would be an empowering step for me since I do not come from a family of academics. My father wants me to become a businessman just like him, possibly taking over his house supplies store. Growing up, I was always impressed by his intuition. While my father and I might not agree on many things, we are both fascinated by risk and strategy. At Harvard, I look forward to refining tools that evaluate uncertainty and democratize access to information. I am eager to contribute to a community that values diversity, academic freedom, and experimental thinking in its teaching and the technologies it pioneers.