Large Scale Learning of Agent Rationality in Two-Player Carnegie Zero-Sum Games Mellon

Chun Kai Ling¹, Fei Fang², J. Zico Kolter^{1,3} Department of Computer Science¹, Institute for Software Research², Carnegie Mellon University Bosch Center for Artificial Intelligence³ chunkail@cs.cmu.edu, feifang@cmu.edu, zkolter@cs.cmu.edu

1. Introduction

- Game theory finds optimal strategies based on known payoffs. Our setting, sometimes known as inverse game theory (Kuleshov, Waugh et al, 2011) is the reverse.
- Our objective is to learn underlying game parameters by observing player actions only.
- Learning utilities allows us to better understand the problem, as opposed to directly predicting strategies.

2. End-to-end learning

- Based on recent work by Ling et. al. (2018)
- Assumes that players act according to the logit Quantal Response Equilibrium (QRE, McKelvey, 1993) of the reduced normal form. • By differentiating `through' game solutions, training is performed end-to-end using SGD to minimize log-loss. • Scales up to larger extensive form games by exploiting the sequence form representation (Von Stengel, 1996).

• Example: Given a context (features) x, learn a payoff matrix P(x)which has equilibrium distributions similar to the actions observed

University





3a. Nested Logit (prior work)

- Logit model: action probabilities based on softmax $u_a^* = \frac{\exp(U_a/\lambda)}{\sum_{a' \in A} \exp(U_{a'}/\lambda)} \quad \text{or} \quad u_a^* = \underset{u \in \Delta_u}{\arg\max u^T U} - \lambda H(u)$
- λ corresponds to the level of (ir)rationality.
- Using logits to model player actions may lead to pathologies (Debreu 1960).
- Nested Logits generalize logits by grouping actions into nests, with different λ for each nest.



- Actions
- Relies on the strong assumption that players play according to the QRE of the reduced normal form.
- Scales poorly in size of the game, requires solving a regularized min-max problem and a linear system for every training point.

3b. Nested Logit QRE

- Obtained by extending nested logits to the multiplayer setting, or alternatively, adding varying levels of rationality at each infoset.
- Natural extension of dilated entropy regularization (Kroer et. al., 2018)

$$\begin{split} \min_{u} \max_{v} u^{T} P v + \sum_{h \in \mathcal{I}_{u}} \lambda_{h} \sum_{a \in \mathcal{A}_{h}} u_{a} \log \frac{u_{a}}{u_{p_{h}}} - \sum_{h \in \mathcal{I}_{v}} \lambda_{h} \sum_{a \in \mathcal{A}_{h}} v_{a} \log \frac{v_{a}}{v_{p_{h}}} \\ \text{subject to} \quad Eu = e, \quad Fv = f. \end{split}$$

u, v: sequence form probs , p_h : action preceding infoset h.

• From a machine learning standpoint: additional λ parameters to be learned. Gradients of loss wrt λ are similar to Ling et. al., which is derived using the implicit function theorem.

$$\nabla_P L = y_u v^T + u y_v^T$$
 $\nabla_{\lambda_h} L = \kappa_h^T y_u$ where

$$(\kappa_{h})_{a} = \begin{cases} 1 + \log(u_{a}/u_{p_{\rho_{a}}}), & \rho_{a} = h \\ -1, & h \in C_{a} \end{cases}$$

$$\begin{pmatrix} \lambda_{\rho_{a}} + \sum_{h' \in \mathcal{C}_{a}} \lambda_{h'} & a = h \\ & & \Gamma u_{n} \rceil = \Gamma - \Xi(u) = P \quad E \quad 0 \rceil^{-1} \Gamma - \nabla_{u} L \rceil$$

Logit

Logit (effectively)

Elimination by aspects (Tversky, 1972)



and quantities involving v are defined analogously to u.

4. Efficient Forward and backward passes

• Utilize first order methods (FOM) of Chambolle and Pock . (2015) to solve convex-concave problems of the form

 $\min_{Ex=x_0} \max_{Fy=y_0} x^T P y + \mathcal{E}(x) - \mathcal{F}(y)$

Requires fast computation of "best response" subproblems

 $\underset{Ex=x_0}{\operatorname{arg\,min}} x^T c_x + \mathcal{E}(x) \quad \text{and} \quad \underset{Fy=y_0}{\operatorname{arg\,min}} y^T c_y + \mathcal{F}(y)$

- Forward pass: subproblem may be solved easily by a single tree traversal or sparse matrix-vector multiplication (Kroer et. al, Hoda et. al., 2010)
- Backwards pass involves solving a large linear system to obtain y. (See 3b).
- Surprisingly, this may be recast to another convex-concave problem by observing that the linear system is precisely the KKT conditions of

5. Experiments

5.1 Synthetic Payoff Matrices

Extensive form games with depth 2 random payoff matrices and \hat{n} actions per player per round.



• Our method is orders of magnitude faster than our previous work

5.2 Information gathering dataset

- The information gathering game (Hunt et. al., 2016) is a one-player 4-stage game which requires players to trade-off between paying for exploration or taking a potentially risky action.
- We used age and education as features and learned



λ -parameters for each of the 4 stages of the game as a function



• The better educated and middle aged demographics enjoy higher rationality (lower λ) corroborating with Hunt et. al...