SIMPLER IS (SOMETIMES) BETTER

A COMPARISON OF COST REDUCING AGENT
ARCHITECTURES IN A SIMULATED BEHAVIORALLY DRIVEN
MULTI-ECHELON SUPPLY CHAIN

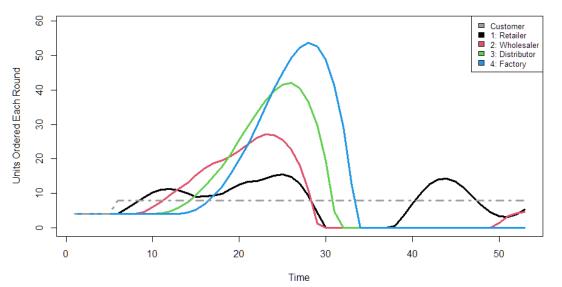


2022 INFORMS Annual Meeting SD34 Supply Chains October 16, 2022

> James Paine jpaine@mit.edu http://jpaine.mit.edu

Bullwhip

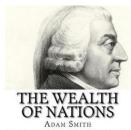




Bullwhip is structurally induced but behaviorally amplified in this environment



Operations Management





ADAM OMITH

Operations Management focus on

18th Century - Division c understanding and improving real

world problems



20th Century – Ford's Assembly Line



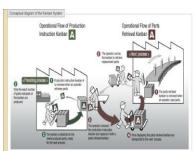
20th Century – Computing Power



20th Century – WWII

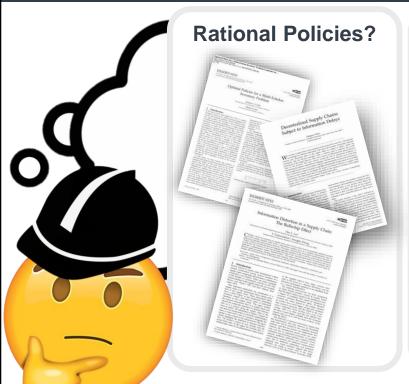


th Century – Taylor's ientific Management'



20th Century – TPS / Lean

What to do?







What to do?

Neural Networks in Supply Chain Management

Horris C. Leung



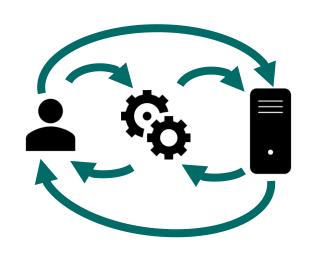
Blackbox !?



This Work



- Assess contribution of different approaches
 - Design agent to follow policies in the context of a well-established simulated supply chain
 - Compare performance outcomes from different policy features







Design of Experiment

| (1) Rule | (2) Degree of | (3) Incentive | (4) Information |
|---------------------|-----------------------|--------------------|-----------------|
| Complexity | Adaptability | Structure | Availability |
| (.1) Base Stock | (.1) Single-Shot | (.1) Self (Myopic) | (.1) Minimal |
| Ordering | . , _ | | , <i>,</i> |
| (.2) Behavioral | (.2) Model-Predictive | (.2) Team (Non- | (2) Low |
| Ordering | Learning | Myopic) | (.2) Low |
| (.3) Deep-Q Network | | | (.3) Standard |
| | | | (.4) High |



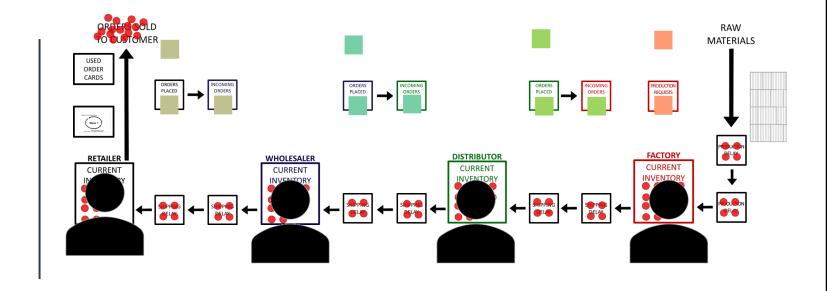
Full Conditions

| Run Code | Agent Feature Description |
|-------------------|--|
| 1.1 - 2.1 | Base Stock Ordering, Single-Shot |
| 1.2 - 2.1 | Behavioral Ordering, Single-Shot |
| 1.3 - 2.1 | Deep-Q Network, Single-Shot |
| 1.3 - 2.2 | Deep-Q Network, Model-Predictive Learning |
| 1.1 - 2.2 - 3.1 - | Base Stock Ordering, Model-Predictive |
| 4.1 | Learning, Self (Myopic), Minimal |
| 1.1 - 2.2 - 3.1 - | Base Stock Ordering, Model-Predictive |
| 4.2 | Learning, Self (Myopic), Low |
| 1.1 - 2.2 - 3.1 - | Base Stock Ordering, Model-Predictive |
| 4.3 | Learning, Self (Myopic), Standard |
| 1.1 - 2.2 - 3.1 - | Base Stock Ordering, Model-Predictive |
| 4.4 | Learning, Self (Myopic), High |
| 1.1 - 2.2 - 3.2 - | Base Stock Ordering, Model-Predictive |
| 4.1 | Learning, Team (Non-Myopic), Minimal |
| 1.1 - 2.2 - 3.2 - | Base Stock Ordering, Model-Predictive |
| 4 2 | Learning, Team (Non-Myopic), Low |

| Run Code | Agent Feature Description |
|-----------------------|---|
| 1.1 - 2.2 - 3.2 - 4.3 | Base Stock Ordering, Model-Predictive Learning, Team (Non-Myopic), Standard |
| 1.1 - 2.2 - 3.2 - 4.4 | Base Stock Ordering, Model-Predictive Learning, Team (Non-Myopic), High |
| 1.2 - 2.2 - 3.1 - 4.1 | Behavioral Ordering, Model-Predictive Learning, Self (Myopic), Minimal |
| 1.2 - 2.2 - 3.1 - 4.2 | Behavioral Ordering, Model-Predictive Learning, Self (Myopic), Low |
| 1.2 - 2.2 - 3.1 - 4.3 | Behavioral Ordering, Model-Predictive Learning, Self (Myopic), Standard |
| 1.2 - 2.2 - 3.1 - 4.4 | Behavioral Ordering, Model-Predictive Learning, Self (Myopic), High |
| 1.2 - 2.2 - 3.2 - 4.1 | Behavioral Ordering, Model-Predictive Learning, Team (Non-Myopic), Minimal |
| 1.2 - 2.2 - 3.2 - 4.2 | Behavioral Ordering, Model-Predictive Learning, Team (Non-Myopic), Low |
| 1.2 - 2.2 - 3.2 - 4.3 | Behavioral Ordering, Model-Predictive Learning, Team (Non-Myopic), Standard |
| 1.2 - 2.2 - 3.2 - 4.4 | Behavioral Ordering, Model-Predictive Learning, Team (Non-Myopic), High |

- Full factorial design on Behavioral x Learning x
- Plus two nonlearning conditions
- Plus two DQN conditions

Specific Simulation Framework: The Beer Game





Simulation Framework

Modular Simulation Model of the Beer Game

- Discrete time simulation
- Modularized to allow different ordering rules from prior literature and to allow embedding in popular AI and simulation tools (Al Gym, Keras, PyTorch, etc)
- E.g. Sterman 89 Ordering rule based on four parameters:

$$O_{t} = MAX(0, \widehat{L}_{t} + \alpha_{S}(S' - S_{t} - \beta SL_{t}) + \varepsilon_{t})$$

$$where \widehat{L}_{t} = \theta L_{t} + (1 - \theta)\widehat{L}_{t-1}$$

O = order placed at time t

L^ = smoothed interpolation of the expected outflow of inventory

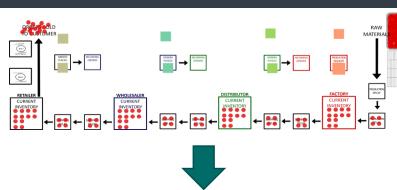
 Θ = smoothing parameter

SL = total inbound supply line of inventory

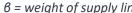
S = current on-hand inventory (or stock)

S' = analogous to the desired or goal on-hand inventory of the player

 θ = weight of supply line



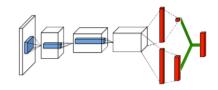
```
Entity Index = 1
or (i in 1:4)
Entity_Index = i
SL = sum(Shipping_flows[Entity_Index..t]) #Total supply line inbound
#L[Entity_Index,t] = Incoming_Order[Entity_Index]
L[Entity_Index,t] = max(min(Incoming_Order[Entity_Index],OH_Inventory[Entity_Index,t]),Order_flows[Entity_Index,1,t-1]
L[Entity_Index,t] = Order_flows[Entity_Index,1,t-1]
  L_hat[Entity_Index.t] = Order_flows[Entity_Index.1.t]
  L_hat[Entity_Index,t-1] = Order_flows[Entity_Index,1,t-1]
  L[Entity_Index,t-1] = L_hat[Entity_Index,t-1]
if (t==20) {
L_hat[Entity_Index,t] = theta[Entity_Index] L[Entity_Index,t] + (1-theta[Entity_Index]) L_hat[Entity_Index,t-1]
S = OH_Inventory[Entity_Index,t] #Stock as of current time
```



Preview of Results

- Incorporating behavioral elements significantly improves performance
- Very complex methods, while interesting, gain performance at the cost of orders of magnitude more complexity
- Agents can be myopic in their optimization goal and have limited information about their environment and still be cost reducing











Single-Shot Policy – Base-Stock Agent





Agent must assume that others are fully rational with same information

- Assuming features of order input are known, costs are increasing along supply chain, and delays are fixed, can follow Clark, A. J., & Scarf, H. (1960)
- This game does not exactly meet that criteria
- Can follow example of Oroojlooyjadid et al '21 (start from C&S then grid search, taking average best values from 50 draws)



Single-Shot Policy – Behavioral Agent



Exogenous $O_{t,c}$ $O_{t,1} = MAX[0, f(\theta_1) + \alpha_1(S_1' - S_{t,1} - \beta_1 SL_{t,1})]$ $O_{t,2} = MAX[0, f(\theta_2) + \alpha_2(S_2' - S_{t,2} - \beta_3)]$

For some entity index i

- Fix $\theta_{n\neq i}$, $\alpha_{n\neq i}$, $S'_{n\neq i}$, $\beta_{n\neq i}$
- Default average performance of Sterman '89 as starting point
- Vary θ_i , α_i , S_i' , β_i to minimize cost function

function
$$O_{t,2} = MAX[0, f(\theta_2) + \alpha_2(S'_2 - S_{t,2} - \beta_3 SL_{t,2})]$$

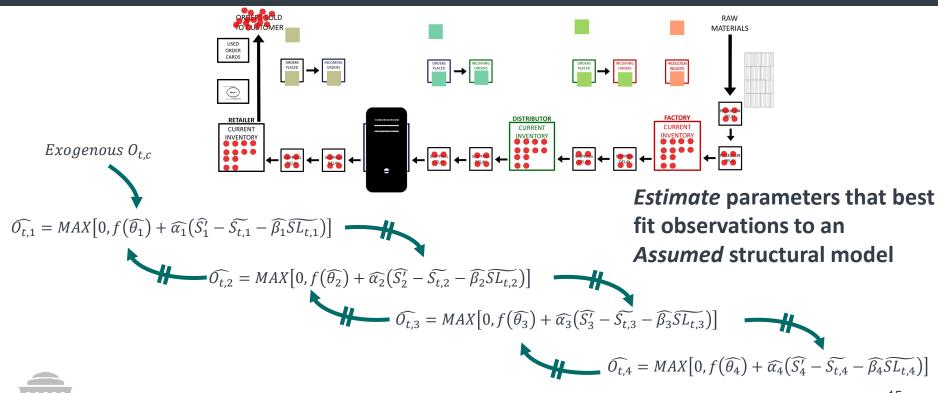
$$O_{t,3} = MAX[0, f(\theta_3) + \alpha_3(S'_3 - S_{t,3} - \beta_3 SL_{t,3})]$$

$$O_{t,4} = MAX[0, f(\theta_4) + \alpha_4(S'_4 - S_{t,4} - \beta_4 SL_{t,4})]$$

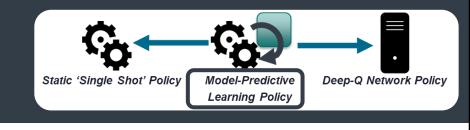


Model-Predictive Learning Policy





Model-Predictive Learning Policy Pseudo Algorithm



Assume system model structure

Calibrate to

Optimize based

and estimated

on assumed

observed

history

t = 0

Assume Structural and Dynamic Model of System

Define Agent position in System model

Define observable space for Agent

Populate initial assumption of parameterization and initializations Define backward calibration memory and forward optimization horizon

for t in 1:horizon

Calibrate System Model given history

ArgMin{System Parameter Estimate}

Error (Observed space of simulation of System Model, Actual Observed space)

Return estimated parameters of System Model

Optimize forward given System Model estimate



model

ArgMax {Agent Decision Rule}

Over t:(t+opt horizon): Reward from t:horizon given System Model estimate

Model-Predictive Learning Policy Information States







Minimal (Narrow Self)

- On-hand inventory position of Agent only
- Most recent order
 placed by Agent only

Low (Self)

- On-hand inventory position of Agent only
- Inbound shipments to
 Agent only
- Order history of the Agent only

Standard (Self + Visible Others)

- On-hand inventory position all positions
- Inbound shipments to all positions
- Order history of the Agent only

High (Near Omniscient)

- On-hand inventory position all positions (all visible state variables)
- positions (all censored state variables)
- Inbound shipments to all positions (all material flows)
- Order history of all positions (all information flows)



Deep Q-Network Policy



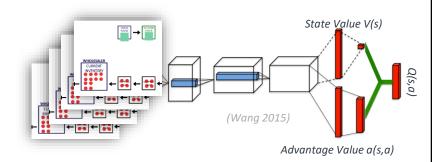


Framework

- OpenAl Gym custom environment based on same functional-form of Beer Game developed in Model-based optimization
- Environment consisting of:
 - All supply chain positions in parallel (versus transfer learning in sequence)
 - Randomly drawn models of real teams (from Sterman '89)
 - Random but bounded simulation horizons
 - Noisy realizations of order decisions

DQN Architecture

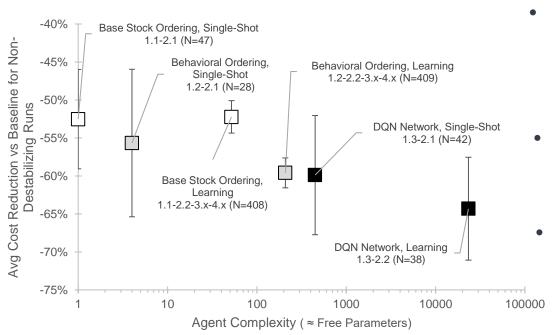
- Dual DQN network (split Action / State Q values)
- Three sequential dense layers with ReLu activation
- Order-plus action space guided by prior model-optimization
- Combination of epsilon-greedy and Boltzmann policy (Wiering 1999)
- Observation space limited to data available in Beer Game (x4 window for sequential memory)







High-Level Results



- Increasing complexity
 yield lower average costs
 (generally) but with
 decreasing returns
 - 'Single-Shot' vs 'Learning' vs 'DQN' offset from each other
 - DQN policies *much* more complex with similar performance



^{*}Average cost reduction at position 2 (wholesaler) vs baseline of no agent with 90% CI bars

^{**}Only cost-reducing results reported

^{***}Information state-varying architectures aggregated

Regressing Predicted Agent Improvement on Policy Features

$$Model~(1): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} = \beta_0 + \beta_1 f_{Behavioral} + \epsilon_{Team_i}$$

Does having a *behavioral* policy matter (single-shot or model-predictive)

$$Model (2): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} = \beta_0 + \beta_1 f_{learning} + \epsilon_{Team_i}$$

Does having a *learning* policy matter (single-shot vs model-predictive)

$$Model~(3): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} = \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{learning} + \epsilon_{Team_i}$$

Does having a behavioral and learning policy matter

$$Model (4): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} = \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{learning} + \beta_3 (f_{Behavioral} * f_{learning}) + \epsilon_{Team_i}$$

Results

$$C_i^{M4} = \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{learning} + \beta_3 (f_{Behavioral} * f_{learning}) + \epsilon_{Team_i}$$

| | Dependent variable: | | | | | |
|-------------------------|---------------------------------------|-----------------------|-------------------------|------------------------|--|--|
| | Performance Improvement over Baseline | | | | | |
| | (1) | (2) | (3) | (4) | | |
| Behavioral Agent | -0.074*** | | -0.071*** | 0.088 | | |
| Learning Agent | | -0.072** | -0.060* | 0.003 | | |
| Behavioral x Learning | | | | -0.173*** | | |
| Constant | -0.522*** | -0.494*** | -0.469*** | -0.525*** | | |
| Observations | 817 | 817 | 817 | 817 | | |
| Log Likelihood | -32.62 | -38.62 | -30.71 | -27.06 | | |
| McFadden R ² | 0.210 | 0.064 | 0.256 | 0.324 | | |
| Residual Std. Error | 0.252 (df = 815) | 0.254 (df = 815) | 0.252 (df = 814) | 0.251 (df = 813) | | |
| | $17.436^{***} (df = 1;$ | $5.291^{**} (df = 1;$ | $10.651^{***} (df = 2;$ | $9.588^{***} (df = 3;$ | | |

- Base agent (Single-Shot Base-Stock Agent) still surprisingly cost reducing
- assumptions is statistically significantly improving

Adding Behavioral

- Adding learning also improves performance
 - Even more helpful when agent is making rationality assumption

F Statistic

Note:

Zooming In: Regressing Predicted Agent Improvement on Information Space in Learning Agents

Does having a *Non-Myopic* (team-level) cost reduction goal affect matter?

$$Model (5): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} = \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{NonMyopic} + \epsilon_{Team_i}$$

$$Learning$$

Does having different levels of information about the environment matter?

$$Model (6): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} = \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{LowInfo} + \beta_2 f_{StandardInfo} + \beta_2 f_{HighInfo} + \epsilon_{Team_i}$$

$$= \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{LowInfo} + \beta_2 f_{StandardInfo} + \beta_2 f_{HighInfo} + \epsilon_{Team_i}$$

$$= \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{LowInfo} + \beta_2 f_{StandardInfo} + \beta_2 f_{HighInfo} + \epsilon_{Team_i}$$

Does having different levels of information about the environment while being Non-Myopic matter?

$$Model (7): \frac{Cost_{Agent_{Team_i}} - Cost_{Baseline_{Team_i}}}{Cost_{Baseline_{Team_i}}} \bigg|_{Learning} = \beta_0 + \beta_1 f_{Behavioral} + \beta_2 f_{NonMyopic} + \beta_3 f_{LowInfo} + \beta_4 f_{StandardInfo} + \beta_5 f_{HighInfo} + \epsilon_{Team_i}$$

Results

М

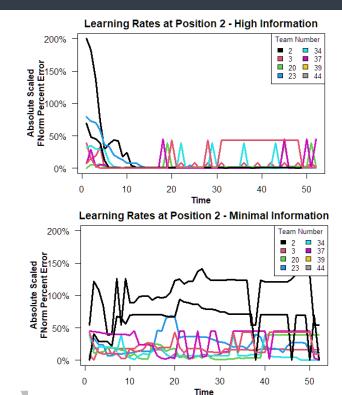
R

$$C_i^{M7} = oldsymbol{eta}_0 + oldsymbol{eta}_1 f_{Behavioral} + oldsymbol{eta}_2 f_{NonMyopic} + oldsymbol{eta}_3 f_{LowInfo} + oldsymbol{eta}_4 f_{StandardInfo} + oldsymbol{eta}_5 f_{HighInfo} + oldsymbol{\epsilon}_{Team_i}$$

Dependent variable:

| | | * | | Again Base learning agent (Bas |
|---|---|--|--|---|
| | Performance Improv (5) | rement over Baseline, g | iven Learning Agent (7) | Stock Learning Myopic Agent) so surprisingly cost reducing |
| Behavioral Agent Non-Myopic Agent | -0.084*** -0.028 | -0.085*** | -0.084*** -0.028 | Adding Behavioral assumptions statistically significantly improvir |
| Low Information Standard Information | 3.62 | -0.016 -0.018 | -0.016 -0.019 | Information availability and team focus not significant |
| High Information Constant | -0.508*** | -0.017 -0.509*** | -0.017 -0.495*** | Having Minimal information comparable to having Omniscient information |
| Observations Log Likelihood | 744 -23.04 | 744 -23.91 | 744 -22.72 | Myopic agents do similar to non- myopic |
| McFadden R ² Residual Std. Error F Statistic | 0.344 $0.250 (df = 741)$ $11.870^{***} (df = 2; 741)$ | 0.313 0.251 (df = 739) 5.481^{***} (df = 4; 739) | 0.345 $0.250 (df = 738)$ $4.865^{***} (df = 5; 738)$ | Note: Results qualitatively similar in most other positions in supply chain. Interaction terms omitted due to lack of significance (see |
| Note: Agent at position 2 (wholesa | aler), Only cost-reducing results | - | 0.1; **p<0.05; ***p<0.01 | due to lack of significance (see paper) |

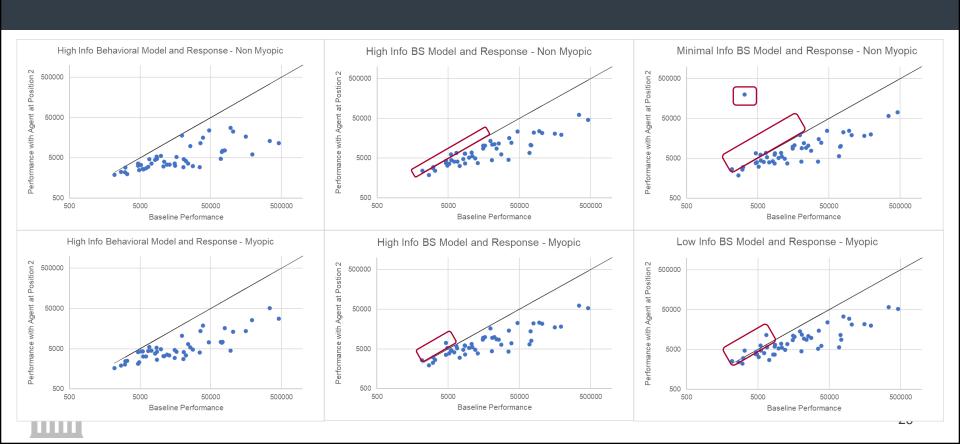
Trade-offs from Low Information



Information Availability Value is Limited

- Key benefit is having a model to fit in the first place
- Dynamically respond to poor initial choices (vs single-shot)
- Near perfect information = more likely to learn 'true' model of environment more quickly
- Minimal information = good enough when entire state can be estimated from limited info

Trade-offs from Low Information



Trade-offs from Low Information

Fraction of Experiments with Destabilizing Outcomes with Model-Predictive Learning Agent in Position 2

| | Information State | | | |
|---|-------------------|-------|----------|------|
| | Minimal | Low | Standard | High |
| $egin{array}{ccc} \mathbf{D} & Base\text{-}Stock \\ \mathbf{D} & (Myopic\ Agent) \end{array}$ | 6.1% | 8.2% | 6.1% | 6.1% |
| Behavioral (Myopic Agent) | 8.2% | 10.2% | 0.0% | 0.0% |

| | Information State | | | |
|-------------------------------|-------------------|-------|----------|------|
| | Minimal | Low | Standard | High |
| Base-Stock (Non-Myopic Agent) | 10.2% | 10.2% | 10.2% | 6.1% |
| Behavioral (Non-Myopic Agent) | 0.0% | 0.0% | 0.0% | 0.0% |

Note:

N = 49 for all conditions

Destabilizing is defined here as -

$$\frac{cost_{Agent_{Team_i}} - cost_{Baseline_{Team_i}}}{cost_{Baseline_{Team_i}}} > 0$$

- More likely in base-stock agents, myopic agents, and with low information
- Did not occur for Behavioral Non-Myopic Agents

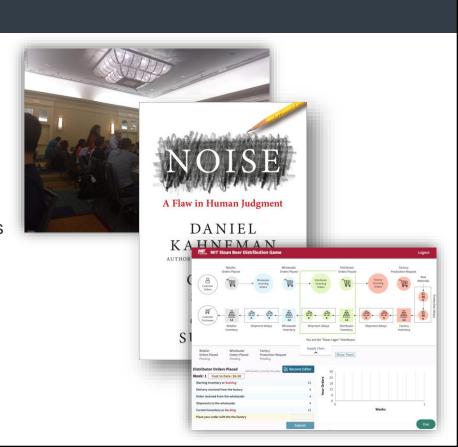




Limitations and Notes

Notable Limitations

- This is an empirically grounded simulation, but not an empirically verified one (yet)
- Simulation of ordering rules from human players, while 'behavioral-based' are still stable
 - Next step is to add noise to ordering realizations
 - Prior experiments imply results should hold, but need to actually test this
- Ordering data from real games implies subtle behavioral differences between 'in-person' and online/hybrid runs.
 - Possibly captured in fitted parameters
 - Future research includes teasing these apart

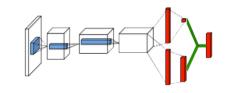


Discussion and Summary

Minor (Methodological) Contributions in this Work

- Developed methods to apply Model-Predictive Learning methods to supply chain ordering policies
- Application of dueling DQN architecture to expand on recent work on DQN in multi-echelon supply chain settings













Main Contributions and Observations of this Work

- Direct comparison of features of policies largely identified separately in prior literature, to reduce cost in behaviorally-driven operational systems
- Incorporating behavioral elements significantly improves performance
- Very complex methods, while interesting, gain performance at the cost of orders of magnitude more complexity
- Agents can be myopic in their optimization goal and have limited information about their environment and still be cost reducing
 - So long as they are *learning*, and assuming behavioral responses low risk of destabilizing



Final Thoughts



Behavioral OM is about Operating Systems that consist of, and exist for, people

- This work shows
 - Cost reducing polices that exist in behaviorally-driven settings
 - How recognizing that people are making decisions is beneficial to these policies
- For managers, this work shows policies that can be locally focused and still achieve globally-preferred outcomes



Thank You!



Please send questions and comment to:

James Paine

jpaine@mit.edu https://jpaine.info

https://github.com/jpain3/ Bullwhip-Policy-Architectures



https://ssrn.com/ abstract=4244188



