

# VC Theory for Inventory Policies

**Yaqi Xie (Chicago Booth)**

joint work with

**Will Ma (Columbia GSB),**

**Linwei Xin (Cornell ORIE)**



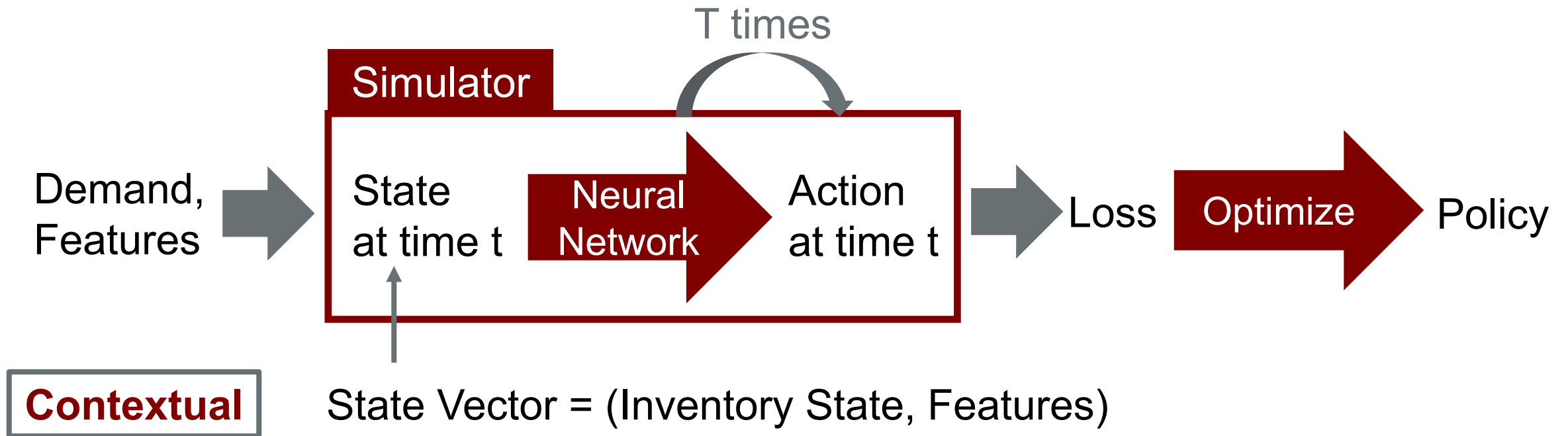
# Background: Data-driven Inventory Control

Approaches	
Classical Approach	Predict and then optimize
Deep Reinforcement Learning (DRL)	End-to-end e.g., <a href="#">Gijsbrechts Boute Van Mieghem Zhang 2022</a>
Supervised Learning (SL)	<a href="#">Sinclair et al. 2023</a> : Hindsight evaluation
	<a href="#">Madeka Torkkola Eisenach Luo Foster Kakade 2022</a> : SL > DRL by Amazon A/B testing
	<a href="#">Alvo Russo Kanoria 2023</a> : SL > heuristics by synthetic data

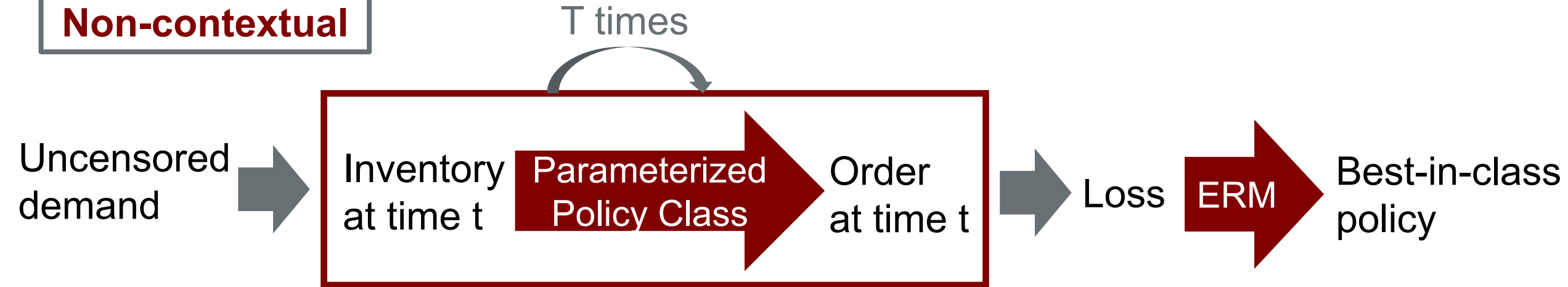
**Our paper: theoretical analysis for classical inventory policies**

# Supervised Learning for Inventory

- Assume uncensored demand  $\leftarrow$  exogenous
- Can evaluate counterfactual performance of any policy on past demand trajectories



## Non-contextual



### Question:

- How badly could we be overfitting? **Estimation/Generalization error**

Summary of “VC Theory for Inventory Policies”: in **SL framework**, we use **VC theory** to analyze the **generalizability** of inventory policies, and how many demand samples are needed to learn **near-optimal** ones.

Sample complexity

# Inventory Basics: Single Durable Good, Periodic Review

- $t = 1, \dots, T$ : finite time horizon
- $x^t$ : inventory at start of time  $t$ , with  $x^1 = 0$
- $y^t$ : inventory after replenishing at time  $t$
- $d^t$ : realized demand at time  $t$
- $c^t$ : cost paid at time  $t$ 
  - $c^t = b \max\{d^t - y^t, 0\} + h \max\{y^t - d^t, 0\} + K \mathbb{I}(y^t > x^t)$
  - $b, h$ : unit costs for understocking, overstocking
  - $K$ : fixed cost for each replenishment

Order quantity

For the talk, we assume **lead time is 0** and **lost-sales** demand:

- $y^t$  can be any number  $\geq x^t$
- $x^{t+1} = \max\{y^t - d^t, 0\}$

Can handle positive lead times for backlogged demand ( $x^{t+1} = y^t - d^t$ )

# Classes of Inventory Policies considered

1)  **$S$  Policies**: defined by **base stock**  $S \in \mathbb{R}_{\geq 0}$

$$y^t = S$$

- optimal if  $d^t$  drawn IID across  $t$  and  $K = 0$

2)  **$(s, S)$  Policies**: defined by  $S \in \mathbb{R}_{\geq 0}$  and **re-order point**  $s \in [0, S]$

$$y^t = x^t + (S - x^t)\mathbb{I}(x^t \leq s)$$

- asymptotically optimal as  $T \rightarrow \infty$  if  $d^t$  drawn IID across  $t$  and  $K > 0$

3)  **$(S^t)$  Policies**: defined by a vector of base stocks  $(S^t)_{t=1}^T \in \mathbb{R}_{\geq 0}^T$

$$y^t = \max\{S^t, x^t\}$$

- optimal if  $d^t$  is independent (non-identical) across  $t$  and  $K = 0$

← Learning  
**near-optimal**  
policy for  
**independent**  
demands

We learn **best-in-class** policies for 1) – 3) **for arbitrary demands**

# Learning Theory Basics

- $\pi \in \Pi$  ( $\Pi = \Pi_S, \Pi_{(S,S)}$ , or  $\Pi_{(S^t)}$ ): inventory policies
- $\mathbf{d} = (d^1, \dots, d^T)$ : demand sequence/trajectory
- $\ell(\pi, \mathbf{d}) = \frac{1}{T} \sum_{t=1}^T c^t(\pi, \mathbf{d})$ : loss of policy  $\pi$  on sequence  $\mathbf{d}$  ← normalized
- $\mathbf{D}$ : (unknown) distribution from which  $\mathbf{d}$  is drawn
  - $L(\pi) = \mathbb{E}_{\mathbf{d} \sim \mathbf{D}}[\ell(\pi, \mathbf{d})]$ : true loss of  $\pi$
  - $\pi^* = \arg \inf_{\pi \in \Pi} L(\pi)$ : policy in  $\Pi$  minimizing true loss ← best-in-class policy
- $(\mathbf{d}_i)_{i=1}^N$ :  $N$  training samples drawn IID from  $\mathbf{D}$ 
  - $\hat{L}(\pi) = \frac{1}{N} \sum_{i=1}^N \ell(\pi, \mathbf{d}_i)$ : empirical loss of  $\pi$
  - $\hat{\pi} = \arg \inf_{\pi \in \Pi} \hat{L}(\pi)$ : policy in  $\Pi$  minimizing in-sample loss ← empirical risk minimizer (ERM)

# Estimation and Generalization Errors

$$\begin{array}{c} \text{Expectation} \\ \text{w.r.t. samples} \end{array} \rightarrow \mathbb{E}[L(\hat{\pi})] = \underbrace{\mathbb{E}[L(\hat{\pi})]}_{\text{Out-of-sample Loss}} = \underbrace{\mathbb{E}[L(\hat{\pi}) - L(\pi^*)]}_{\text{ERM Best-in-class Estimation Error}} + L(\pi^*) \leftarrow \begin{array}{c} \text{Not depend} \\ \text{on samples} \end{array}$$

$$\underbrace{\mathbb{E}[L(\hat{\pi}) - L(\pi^*)]}_{\text{Estimation Error}} \leq \underbrace{\mathbb{E} \left[ \sup_{\pi \in \Pi} (L(\pi) - \hat{L}(\pi)) \right]}_{\text{Generalization Error}} = o \left( \sqrt{\text{PDim}(\Pi)/N} \right) \quad \forall \text{ distributions } \mathbf{D}$$

Goal: bound the **generalization errors** of  $S$ ,  $(s, S)$  and  $(S^t)$  policy classes

- w.r.t. sample size  $N$  and horizon length  $T$
- **without any assumption** on the demand distribution  $\mathbf{D}$

# Generalization Error, VC/Pseudo-dimension

VC Theory:  $\underbrace{\mathbb{E}[L(\hat{\pi}) - L(\pi^*)]}_{\text{Estimation Error}} \leq \underbrace{\mathbb{E} \left[ \sup_{\pi \in \Pi} (L(\pi) - \hat{L}(\pi)) \right]}_{\text{Generalization Error (GE)}} = o \left( \sqrt{\text{PDim}(\Pi)/N} \right) \quad \forall \text{ distributions } \mathbf{D}$

## Our results:

- 1)  $\text{PDim}(\Pi_S) \leq 2$
- 2)  $\text{PDim}(\Pi_{(s,S)}) = O(\log T)$ ,  $\text{GE}(\Pi_{(s,S)}) = \Omega(\sqrt{\log T / \log \log T / N})$
- 3)  $\text{PDim}(\Pi_{(S^t)}) = \Omega(T)$ ,  $\text{GE}(\Pi_{(S^t)}) = O(\sqrt{1/N})$

	Policy Class	Lower Bound	Upper Bound
New results in <b>Red</b>	$S$ Policies	$\Omega(\sqrt{1/N})$	$o(\sqrt{1/N})$
	$(s, S)$ Policies	$\Omega(\sqrt{\log T / \log \log T / N})$	$o(\sqrt{\log T / N})$
	$(S^t)$ Policies	$\Omega(\sqrt{1/N})$	$o(\sqrt{1/N})$

Arbitrary demand

Surprising improvement from literature

Independent demand

# $S$ Policies, $(s, S)$ Policies

Policy Class	Lower Bound	Upper Bound
$S$ Policies	$\Omega(\sqrt{1/N})$	$O(\sqrt{1/N})$
$(s, S)$ Policies	$\Omega(\sqrt{\log T / \log \log T / N})$	$O(\sqrt{\log T / N})$

Novel way to prove  
Newsvendor (uncensored)

←  $\text{PDim}(\Pi_S) \leq 2$

←  $\text{PDim}(\Pi_{(s,S)}) = O(\log T)$


$\exists \text{instance, s.t.}$   
 $\text{PDim}_\gamma(\Pi_{(s,S)}) = \Omega(\log T / \log \log T)$

Fan Zhou 2024 derive  $O(\sqrt{T/N})$   
 for IID integer demands

- Policy class  $\Pi$  **shatters** samples  $\mathbf{d}_1, \dots, \mathbf{d}_m$  with **witnesses**  $\tau_1, \dots, \tau_m$ , if for all  $G \subseteq \{1, \dots, m\}$  (“Good”), there exists  $\pi \in \Pi$  such that
 
$$\begin{aligned} \ell(\pi, \mathbf{d}_i) &\leq \tau_i \quad \forall i \in G; \\ \ell(\pi, \mathbf{d}_i) &> \tau_i \quad \forall i \notin G \end{aligned}$$
- Pseudo-dimension** of  $\Pi$ , or  $\text{PDim}(\Pi)$ , is the maximum of samples that it can shatter with witnesses

$(S^t)$  Policies:  $\text{GE}(\Pi_{(S^t)}) = \mathcal{O}(\sqrt{1/N})$

$$\underbrace{\mathbb{E} \left[ \sup_{\pi \in \Pi} (L(\pi) - \hat{L}(\pi)) \right]}_{\text{Generalization Error (GE)}} \leq \boxed{\text{Rademacher Complexity}} \stackrel{\text{Lipschitz, Talagrand}}{\leq} \boxed{\text{Rademacher for Inventory Level}} \leq \gamma\text{-shattering dimension } \text{PDim}_{\gamma}(Y(\Pi))$$

Ending inventory by  $\Pi$  

### $\gamma$ -shattering dimension:

- $Y(\Pi)$  **shatters** samples  $d_1, \dots, d_m$  with **witnesses**  $\tau_1, \dots, \tau_m$  **at scale**  $\gamma > 0$ , if for all  $G \subseteq \{1, \dots, m\}$  (“Good”), there exists  $\pi \in \Pi$  such that

$$\begin{aligned} \text{Ending inventory} &\begin{cases} y(\pi, d_i) \leq \tau_i - \gamma & \forall i \in G; \\ y(\pi, d_i) > \tau_i + \gamma & \forall i \notin G \end{cases} \end{aligned}$$

- $\text{PDim}_{\gamma}$  is the maximum  $m$  of samples that it can shatter with witnesses **at scale**  $\gamma$

# $\text{PDim}_\gamma \left( \Pi_{(s^t)} \right) = \mathbf{O}(1/\gamma)$ : Proof Sketch I

We prove: max # of trajectories does not grow with  $T$

Necessary  
→

Demand  
trajectories

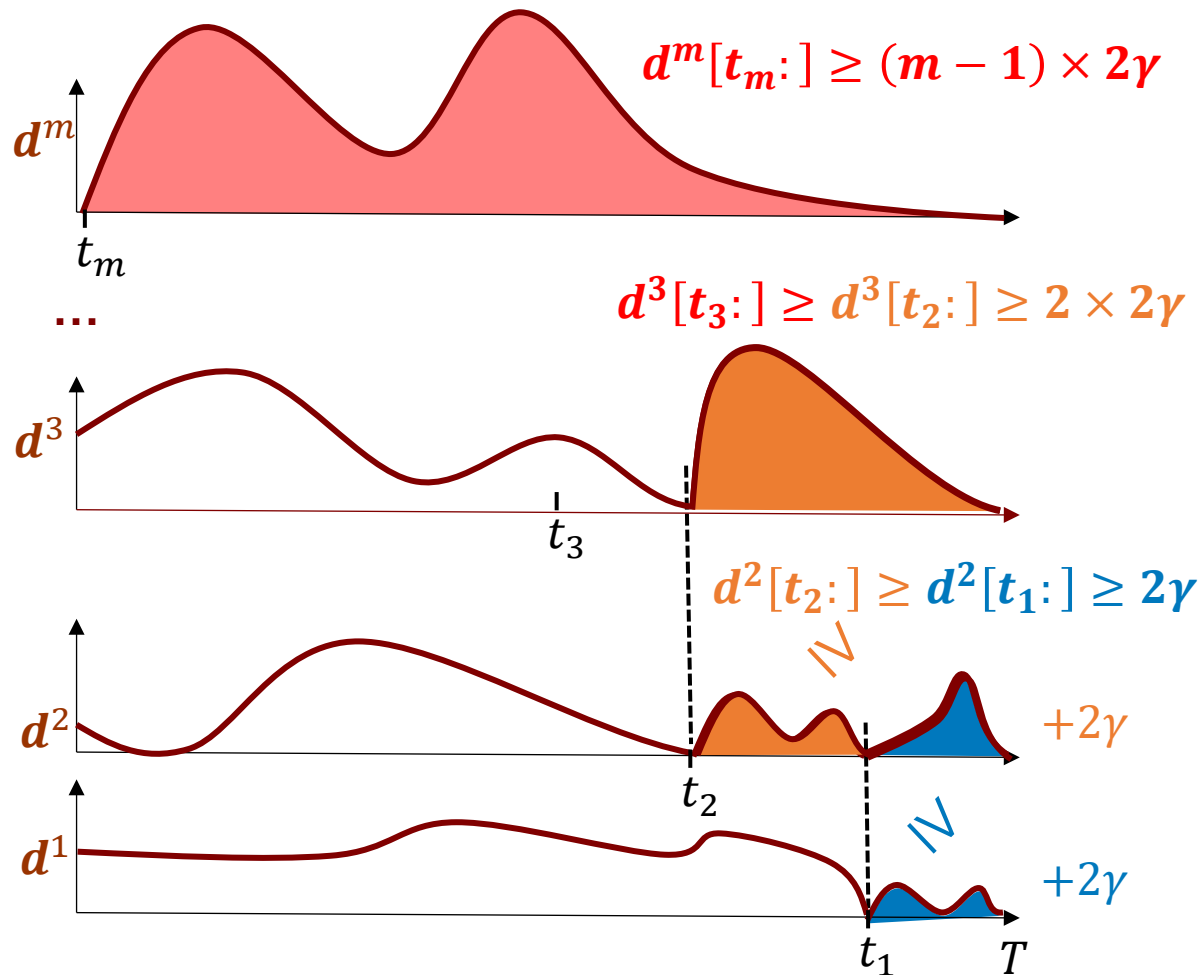
Ending Inventory	policies				
	$\pi^1$	$\pi^2$	$\pi^3$	...	$\pi^m$
$d^1$	$> \tau_1 + \gamma$	$\leq \tau_1 - \gamma$	$\leq \tau_1 - \gamma$		$\leq \tau_1 - \gamma$
$d^2$	$\leq \tau_2 - \gamma$	$> \tau_2 + \gamma$	$\leq \tau_2 - \gamma$		$\leq \tau_2 - \gamma$
$d^3$	$\leq \tau_3 - \gamma$	$\leq \tau_3 - \gamma$	$> \tau_3 + \gamma$		$\leq \tau_3 - \gamma$
...					
$d^m$	$\leq \tau_m - \gamma$	$\leq \tau_m - \gamma$	$\leq \tau_m - \gamma$		$> \tau_m + \gamma$

- Assume  $\tau_1 = \dots = \tau_m =: \tau$  for proof sketch
- Let  $t_i$  be the last replenishment point for  $\pi^i$  on  $d^i$ , and  $d^i[t_i:] := \sum_{t'=t}^T d_{t'}^i$ , be demand from that point onward

$$\pi_{t_i}^i - d^i[t_i:] = y_T(\pi^i, d^i) > \tau + \gamma$$

$$\pi_{t_i}^i - d^j[t_i:] \leq y_T(\pi^i, d^j) \leq \tau - \gamma \quad \forall j \neq i \implies d^j[t_i:] - d^i[t_i:] > 2\gamma \quad \forall j \neq i$$

# $\text{PDim}_\gamma \left( \Pi_{(\text{st})} \right) = O(1/\gamma)$ : Proof Sketch II



Re-index so that  $t_m < \dots < t_1$ ,

$$d^j[t_i:] > d^i[t_i:] + 2\gamma \quad \forall i \neq j \quad \star$$

But  $1 \geq \pi_{t_m}^m \geq d^m[t_m:]$  because  $t_m$  is a last replenishment point on  $d^m$ !

Hence can shatter at most  $m = O(1/\gamma)$  trajectories, for any  $\gamma > 0$ .

# $(S^t)$ Policies: “Horizon-free” Learning Guarantees

	Independent Demands	Arbitrary Demands	
<b>Data-driven inventory:</b> Empirical dynamic program or Genetic RL	Levi Roundy Shmoys 07: $O(\sqrt{T^5 \log T / N})$	Shapiro Dentcheva Ruszczynski 09: $O(\sqrt{T/N})$  Shalev-Shwartz Shamir Srebro Sridharan 10: $O(\sqrt{T/N})$	<b>Stochastic Optimization -</b> SAA for $T$ -dim decision space: Interpolation or Covering number
	Cheung Simchi-Levi 19: $O(\sqrt{T^6 \log T / N})$		
	Qin Simchi-Levi Zhu 23: $O(\sqrt{T/N})$		
<b>Episodic RL:</b> Finite state and action spaces	Yin Bai Wang 21: $O(\sqrt{T \log T / N})$		
	Zhang Ji Du 22: $O(\sqrt{T \log N / N})$		
	<b>Our Paper: <math>O(\sqrt{1/N})</math></b>		
		↑ <b>VC Theory</b>	

## Conclusion

Our paper provides the theoretical analysis for the **supervised learning** framework:

we use **VC theory** to analyze the **generalization error/sample complexity** of **inventory** policies.

Policy Class	Lower Bound	Upper Bound
$S$ Policies	$\Omega(\sqrt{1/N})$	$O(\sqrt{1/N})$
$(s, S)$ Policies	$\Omega(\sqrt{\log T / \log \log T / N})$	$O(\sqrt{\log T / N})$
$(S^t)$ Policies	$\Omega(\sqrt{1/N})$	$O(\sqrt{1/N})$

Arbitrary demand  
Surprising improvement from literature  
Independent demand

- VC theory is a powerful tool.
- The number of policy parameters may not be the measure of overfitting.

# Thanks for Your Attention!

VC Theory for Inventory Policies

Yaqi Xie, Will Ma, Linwei Xin

[https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=4794903](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=4794903)