

# Domain Transfer for Reinforcement Learning Agents

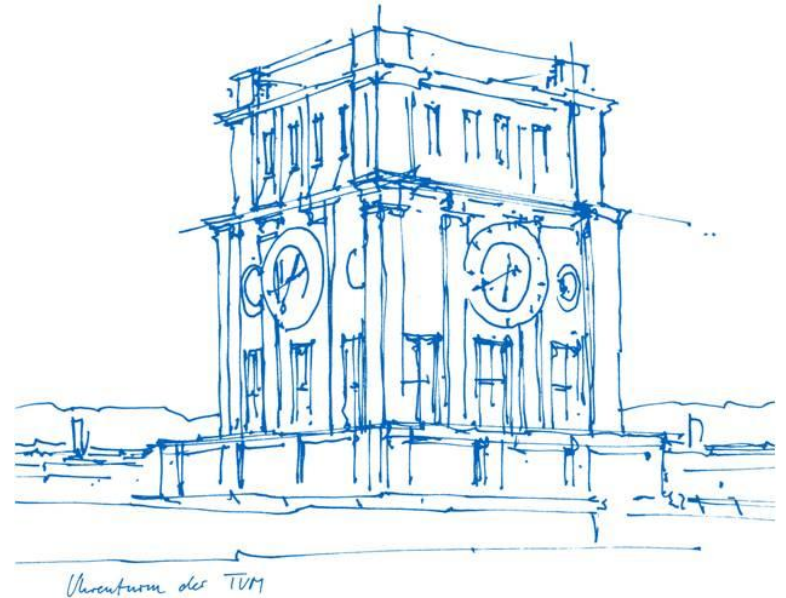
## Final Presentation - PreciBake

Ahmad Qasim, Murat Can Karacabey  
Abdul Moeed, Sebastian Oehme

Mentors: M. Sundholm, H. Belhassan  
TUM-Mentor: M. Rauchensteiner

Data Innovation Lab  
Department of Mathematics  
Technische Universität München

München, 17. February 2020



Spelt Sourdough \$6.50

Honey Spelt Bread \$6.50

5-Grain Pan Bread \$5.50

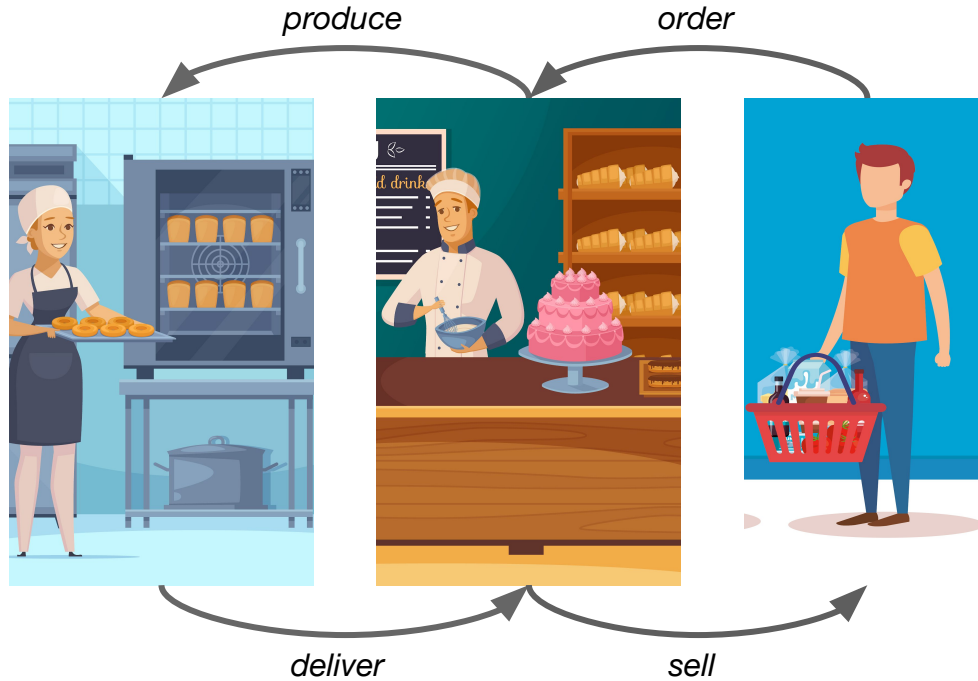


Simple Sourdough \$6.00

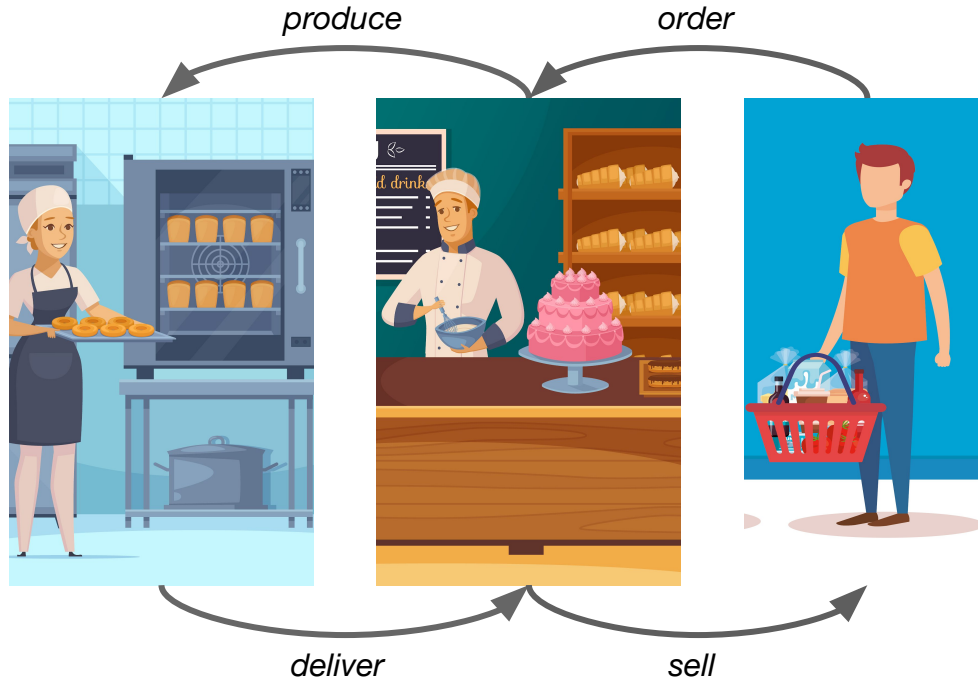
Ancient Multigrain Pan Bread \$6.00



# Inventory Control



# Inventory Control



## Main goal:

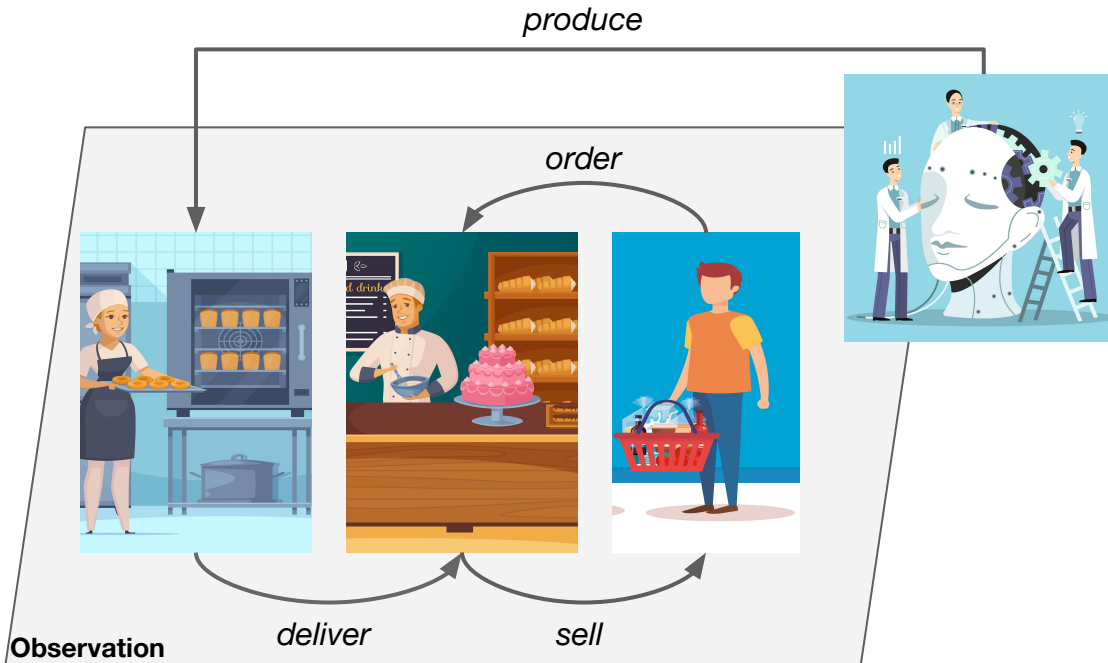
- Maximize product sales

## Sub-goals:

- Minimize waste
- Maximize freshness
- Minimize waiting time for customer



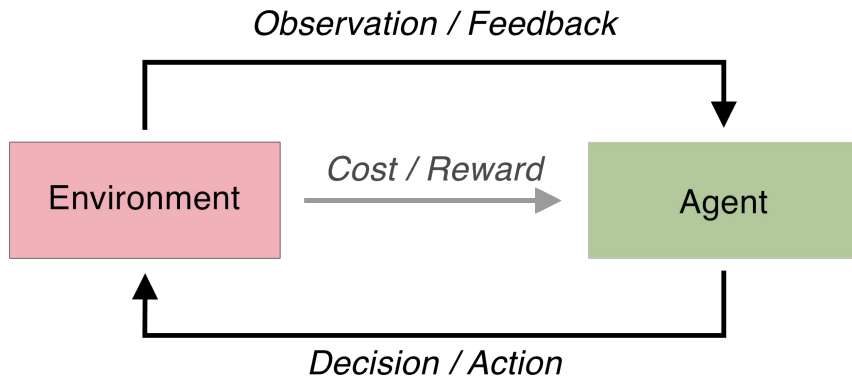
# Inventory Control with Reinforcement Learning



Create an agent that

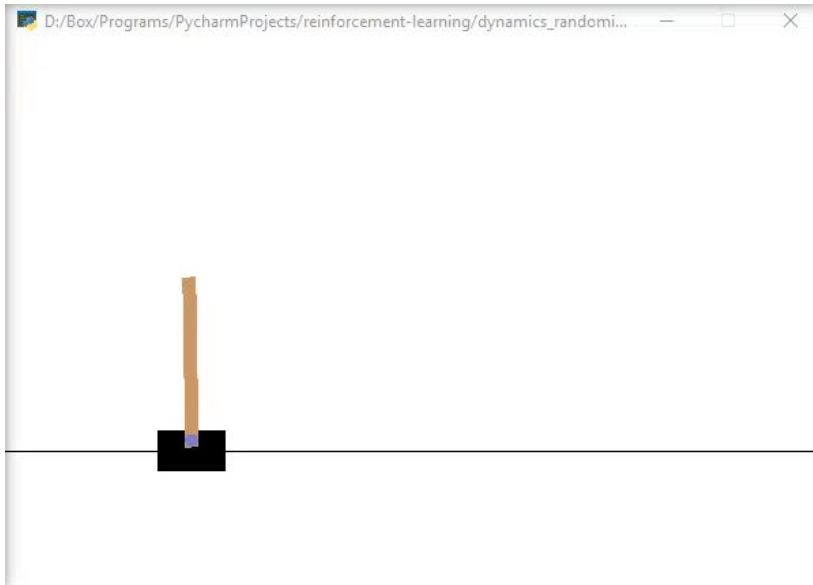
- Satisfies problem constraints
- Safely operates in real world after training on simulation (domain transfer)

# Reinforcement Learning (RL)



On each ***episode***, agent must take ***actions*** that maximize expected ***reward*** for each ***state***

# Example - RL

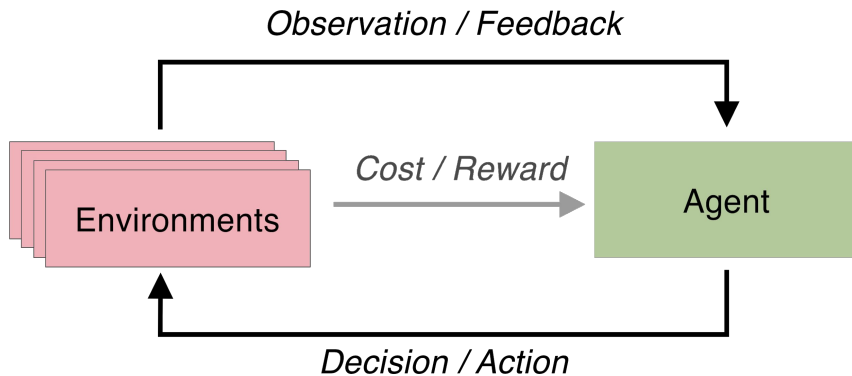


Cartpole - OpenAI Gym

## Task: Balance pole on cart

- **State**
  - Cart position, velocity
  - Pole angle
- **Actions**
  - *Move left or right*
- **Reward**
  - *+1 for each time step*
  - *0 when pole falls*

# RL with Dynamics Randomization

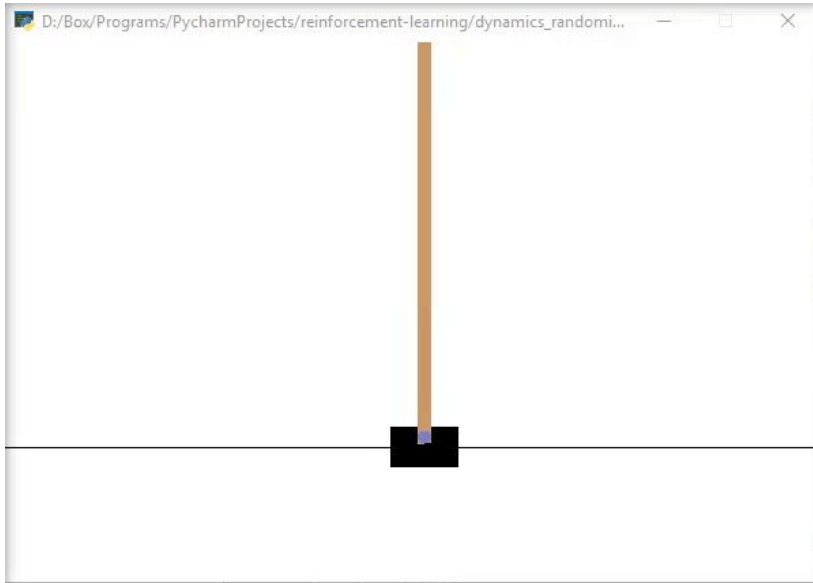


**Idea:** To generalize to new environments, sample a different environment on each episode

**Practice:** Environment is defined by certain parameters, sample different values of these parameters on each episode



# Dynamics Randomization - Example



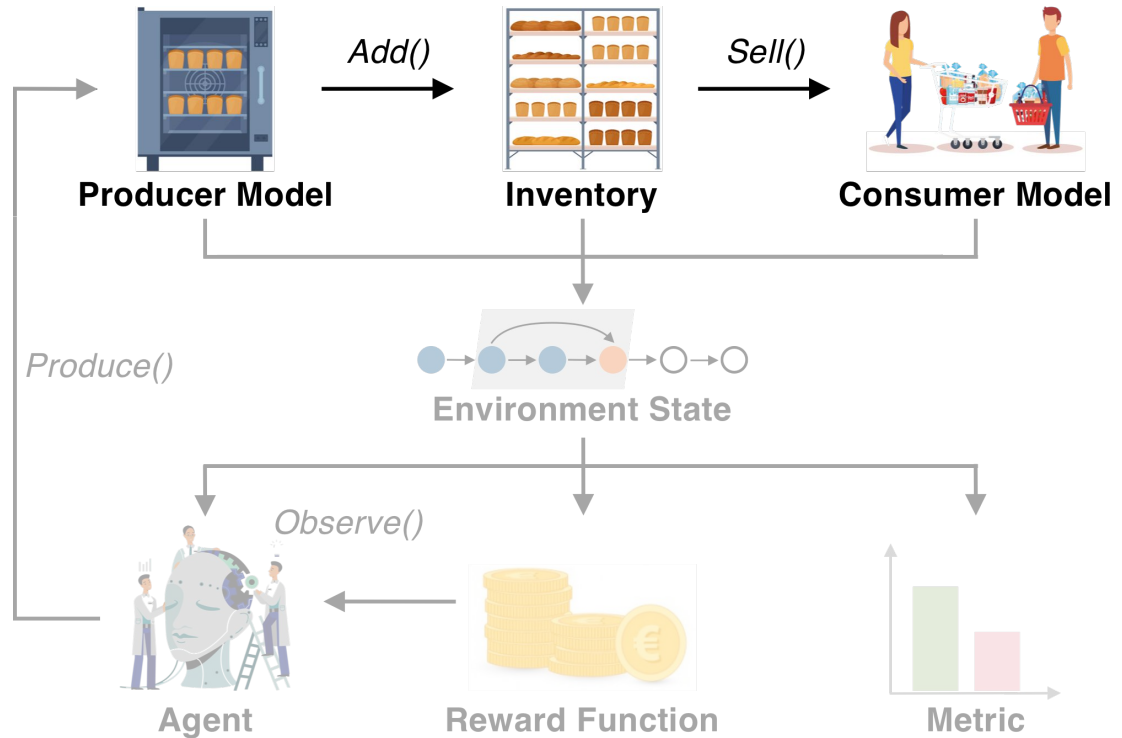
Cartpole - OpenAI Gym

*Per episode, sample*

- **Pole Length:** [0.3, 0.6]
- **Pole Mass:** [0.5, 1.5]

# Environment

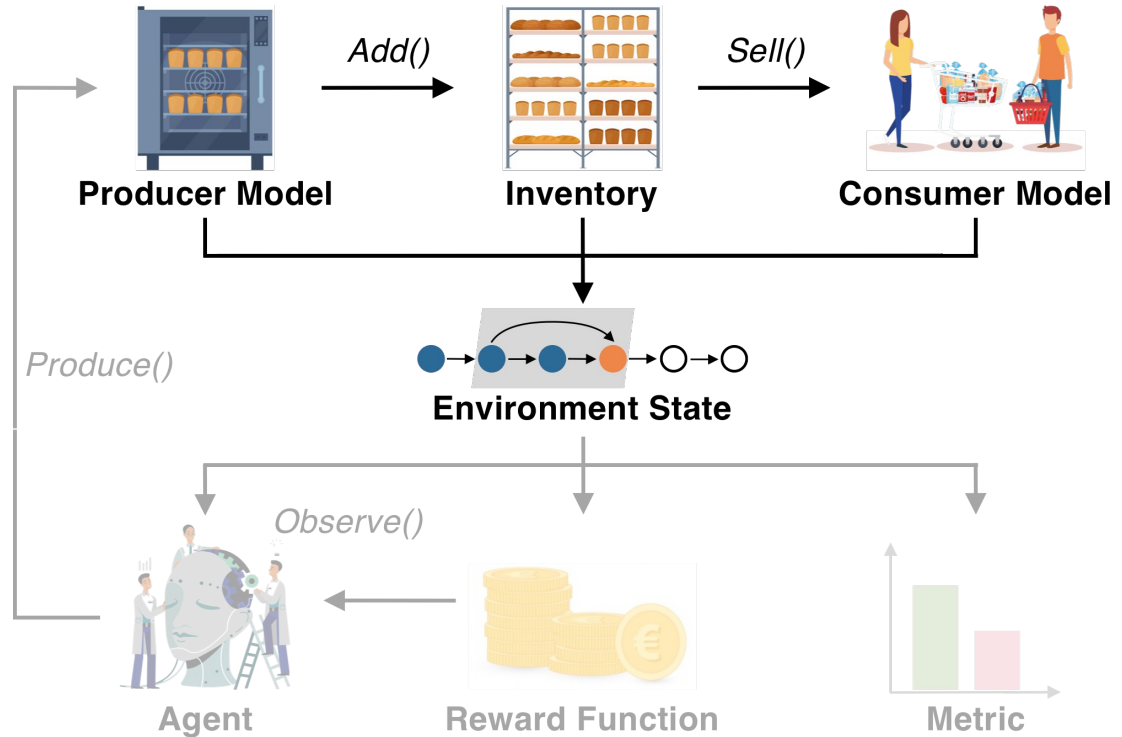
- **Producer Model**
  - Time-shift
  - Quantity limit
- **Inventory**
  - Storage limit
  - Freshness
- **Consumer Model**
  - Stochastic



# Environment

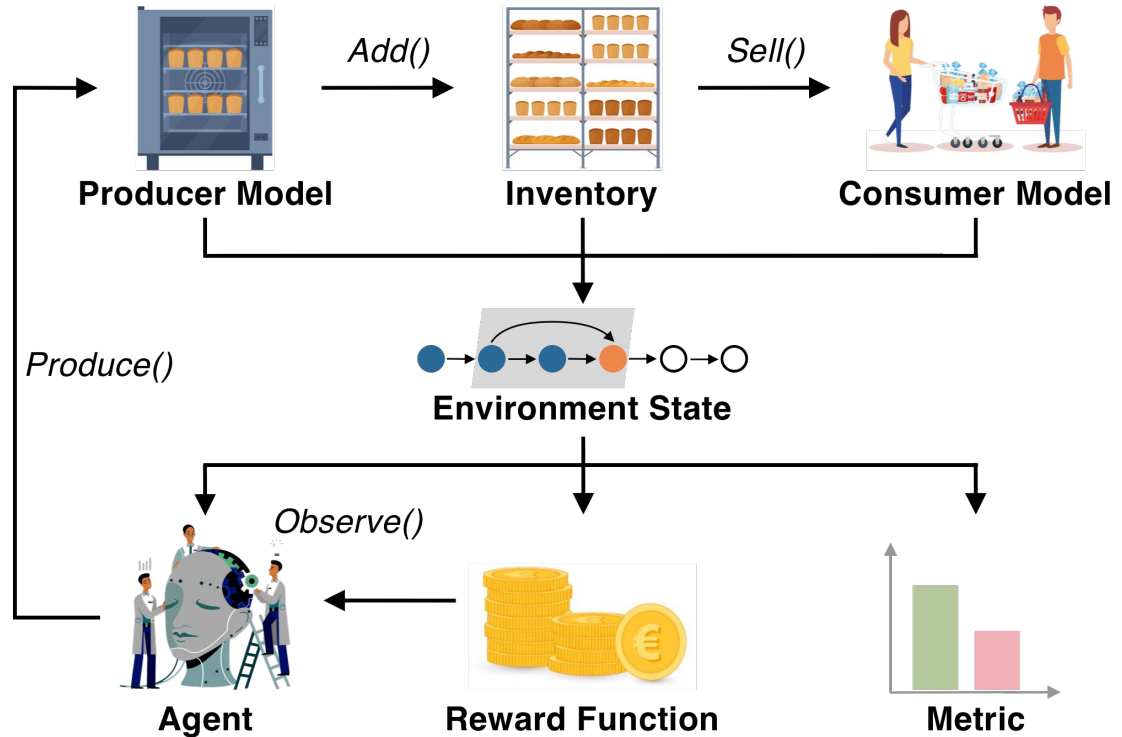
- **Environment State**

- Production queue
- Order queue
- Inventory
  - Number of the items
  - Age of the items



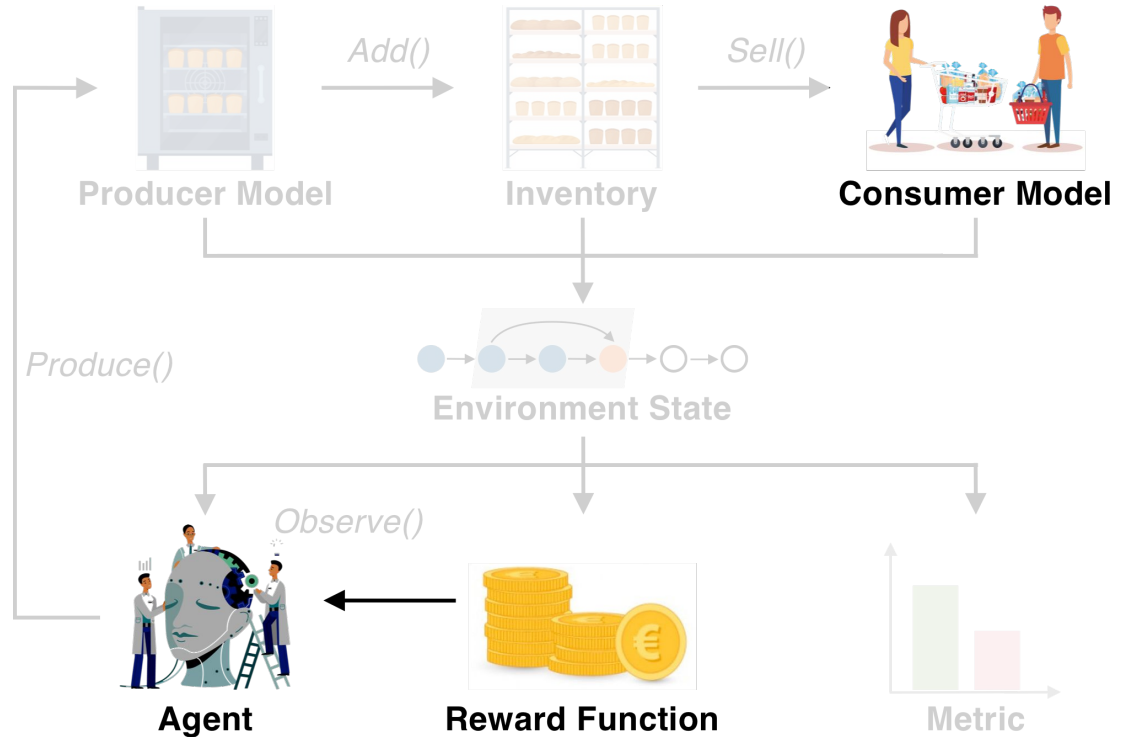
# Environment

- **Agent**
  - Decides on what to produce
- **Reward**
  - Agent-specific evaluation of the state
- **Metric**
  - Common evaluation of the agents' performances



# Environment - Our Task

- **Main task:**
  - Create a robust agent that can perform well in varying demand trends
- **Requires:**
  - An expressive consumer model



# Consumer Model

- Simulate behavior of giving orders
- Requirements for the model
  - Expressibility
  - Flexibility
- Problem
  - No real world data available



# Consumer Model

- Simulate behavior of giving orders
- Requirements for the model
  - Expressibility
  - Flexibility
- **Problem**
  - **No real world data available**

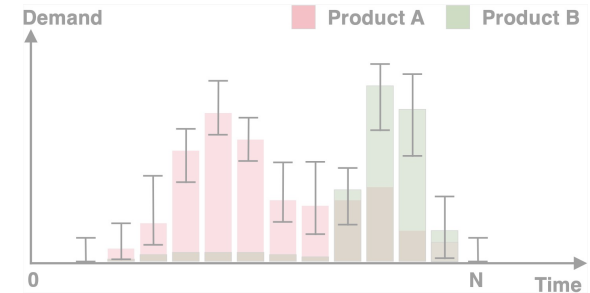
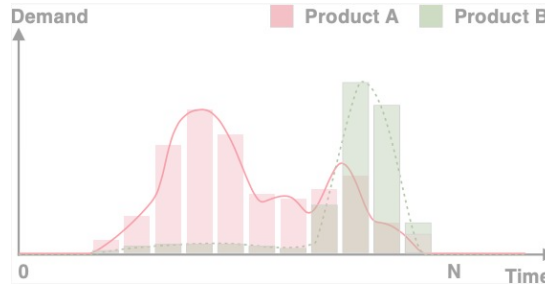
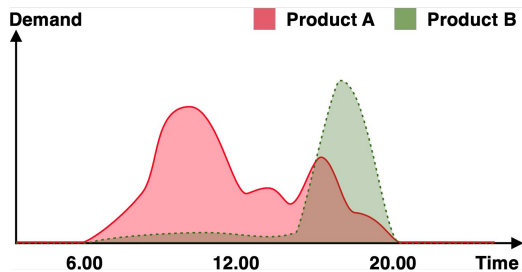


## **Solution approaches**

1. Learn a model from a public dataset
2. Use a generative model to generate data

# Generative Consumer Model

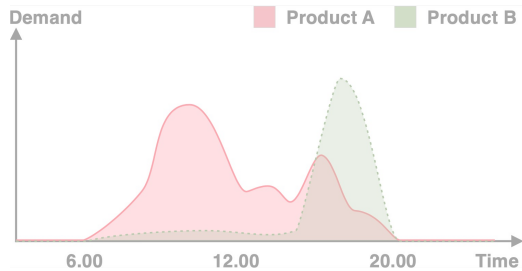
A generative model that is based on Poisson Distribution



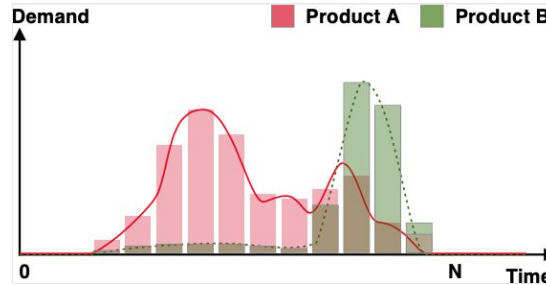
- Example expected demand

# Generative Consumer Model

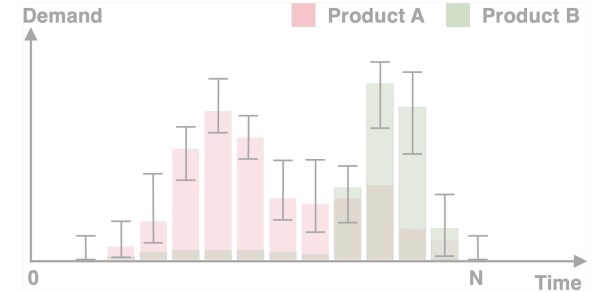
A generative model that is based on Poisson Distribution



- Example expected demand

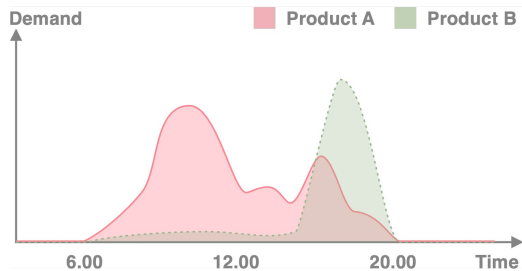


- Create bins with means
- Generate demand

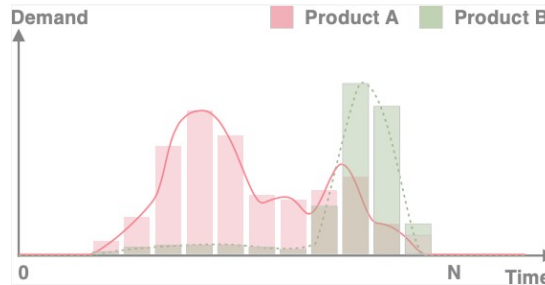


# Generative Consumer Model

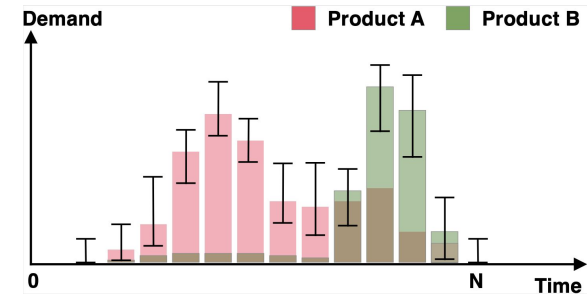
A generative model that is based on Poisson Distribution



- Example expected demand



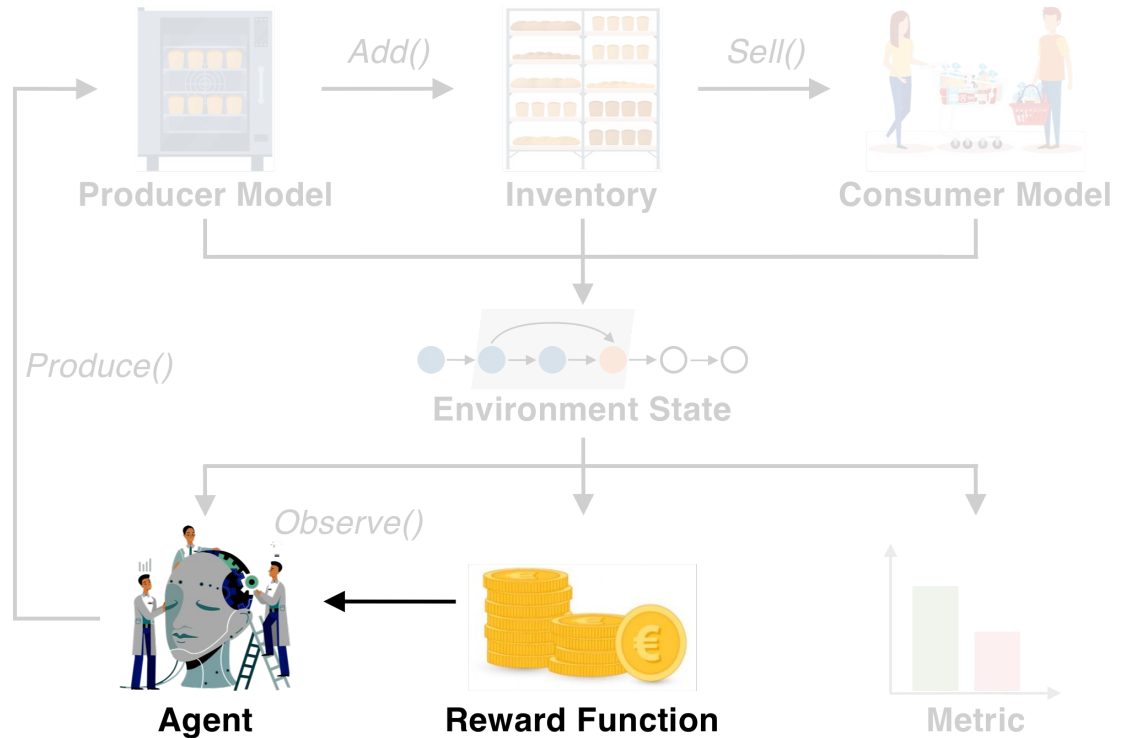
- Create bins with mean
- Generate demand



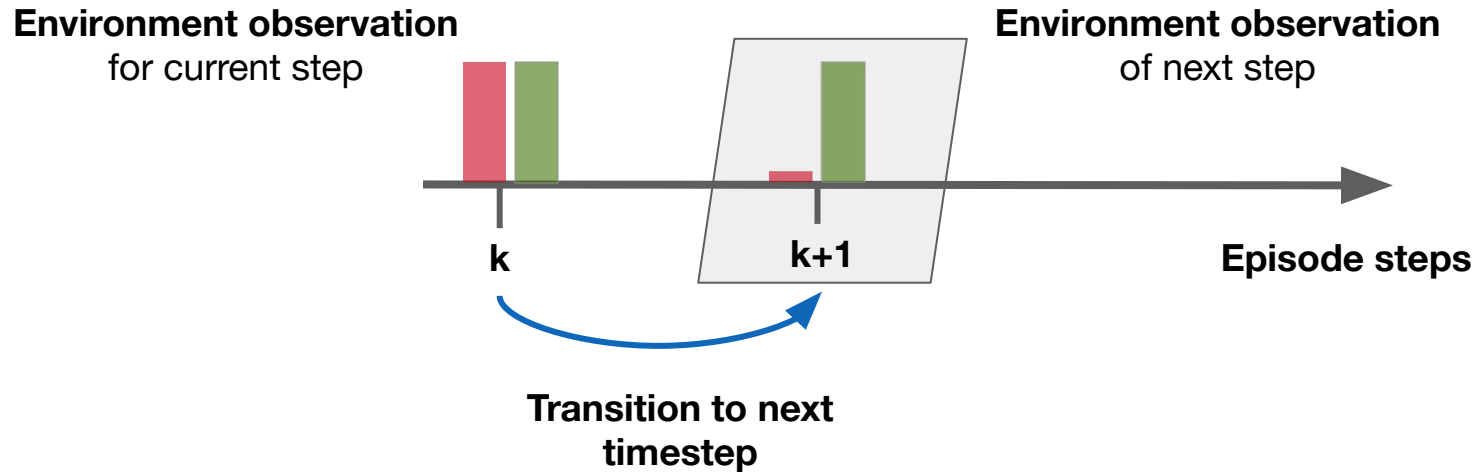
- Vary means for dynamics randomization

# Environment

- **Agent** has to choose
  - Product type
  - Product quantity



# Agents Overview



■ Product A ■ Product B



# Agents

## Baseline

- Benchmark for comparison
- Simulates a real bakery decision process
- Fill inventory if below threshold

## Deep Deterministic Policy Gradient (DDPG)

- Model-free
- High-performing RL algorithm
- Using Deep RL to estimate best action
- Reward engineering

## Dynamic Programming (DP) with Demand Prediction

- Model-based
- Adapted from D. Bertsekas
- Interpretable and follows the principle of optimality
- Cost function

# Agents

## Baseline

- Benchmark for comparison
- Simulates a real bakery decision process
- Fill inventory if below threshold

## Deep Deterministic Policy Gradient (DDPG)

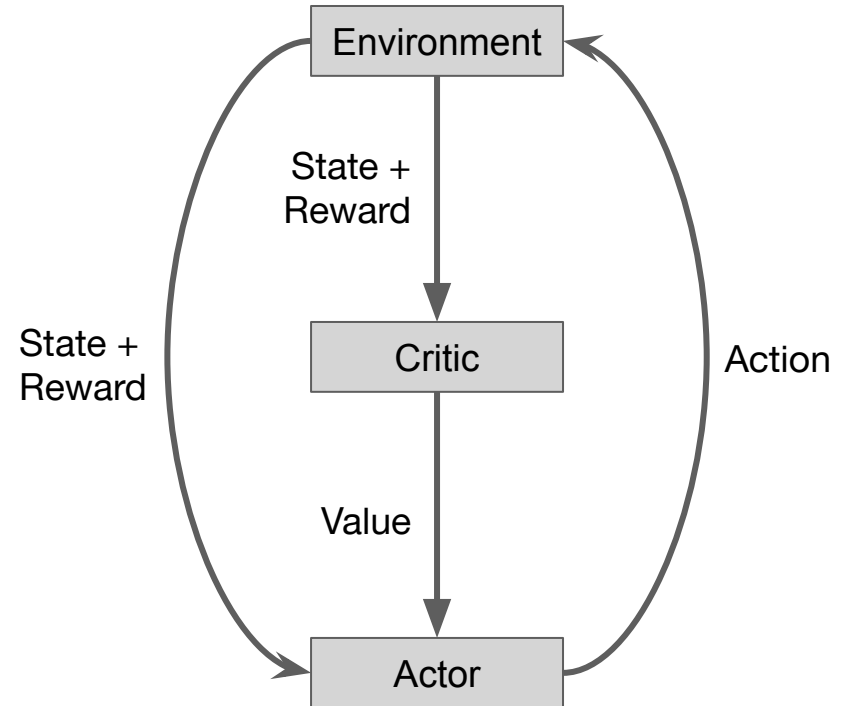
- Model-free
- High-performing RL algorithm
- Using Deep RL to estimate best action
- Reward engineering

## Dynamic Programming (DP) with Demand Prediction

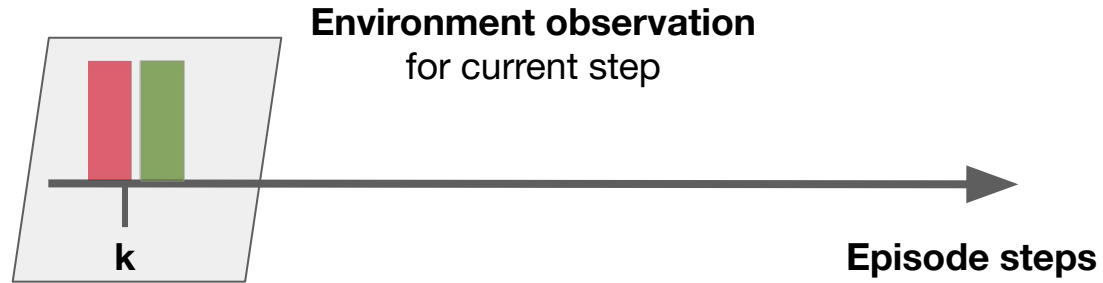
- Model-based
- Adapted from D. Bertsekas
- Interpretable and follows the principle of optimality
- Cost function

# DDPG

- Environment has a state and returns a reward
- Critic estimates the value of state
- Actor generates actions
- Challenges:
  - Several hyper-parameters to tune
  - Reward engineering

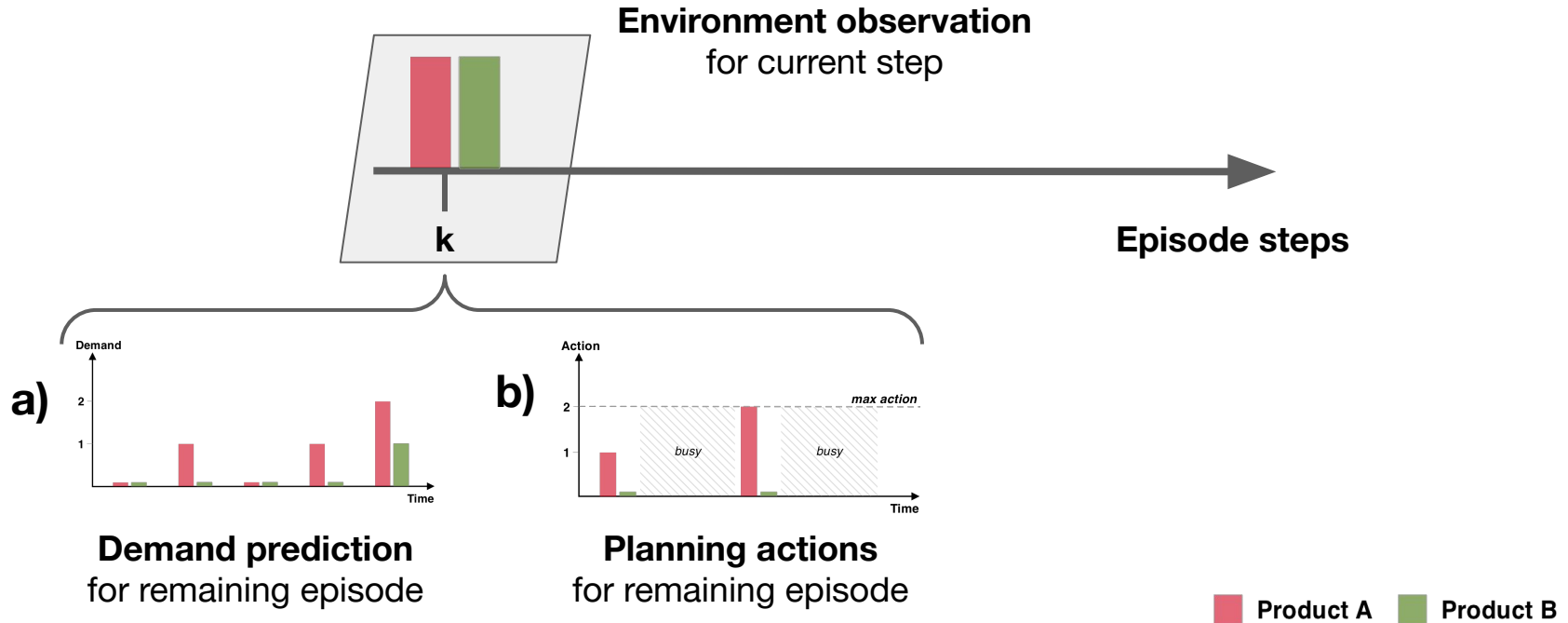


# DP Agent overview



 Product A  Product B

# DP Agent overview



# Demand Model

- Predicts future consumer orders
  - DP agent acts based on it
  - DP agent's performance, highly dependent on demand model
- Tried different approaches



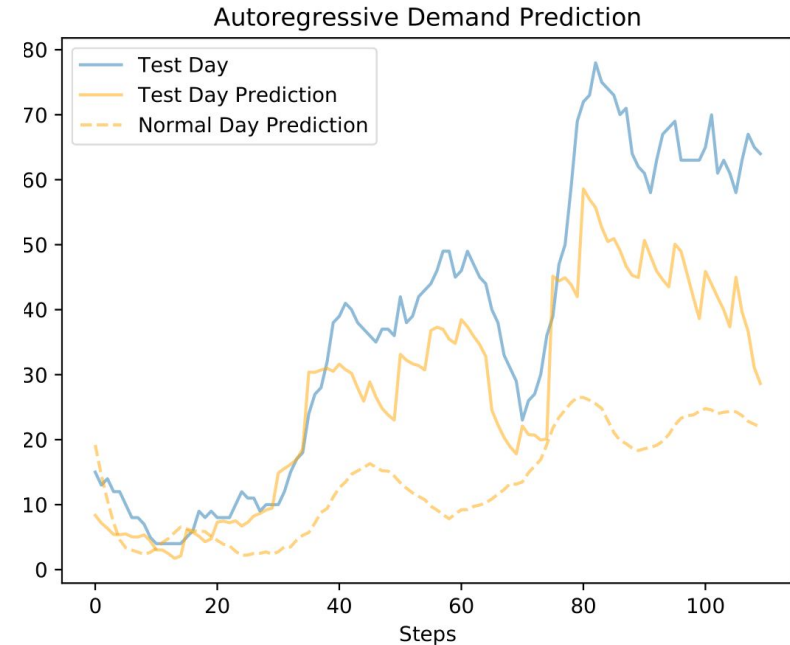
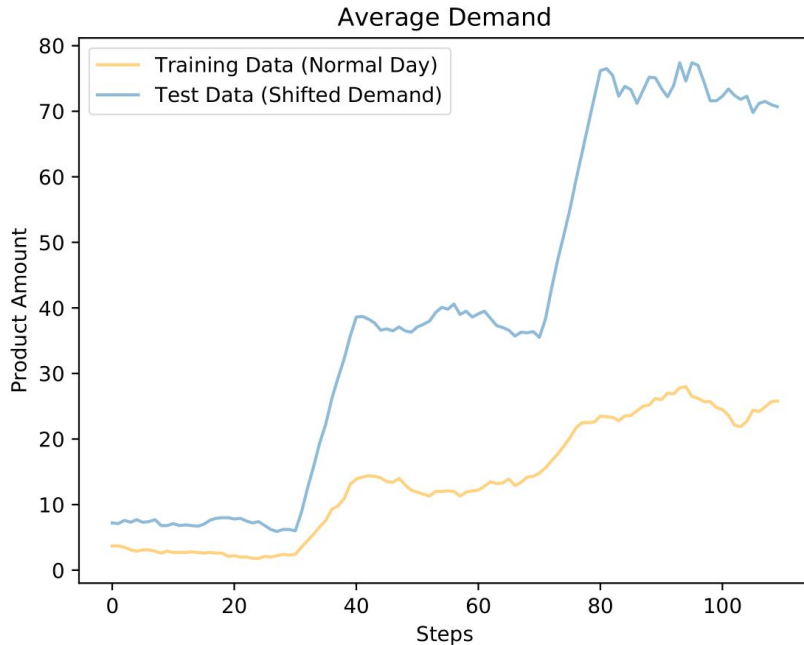
# Demand Model

- Predicts future consumer orders
  - DP agent acts based on it
  - DP agent's performance, highly dependent on demand model
- Tried different approaches

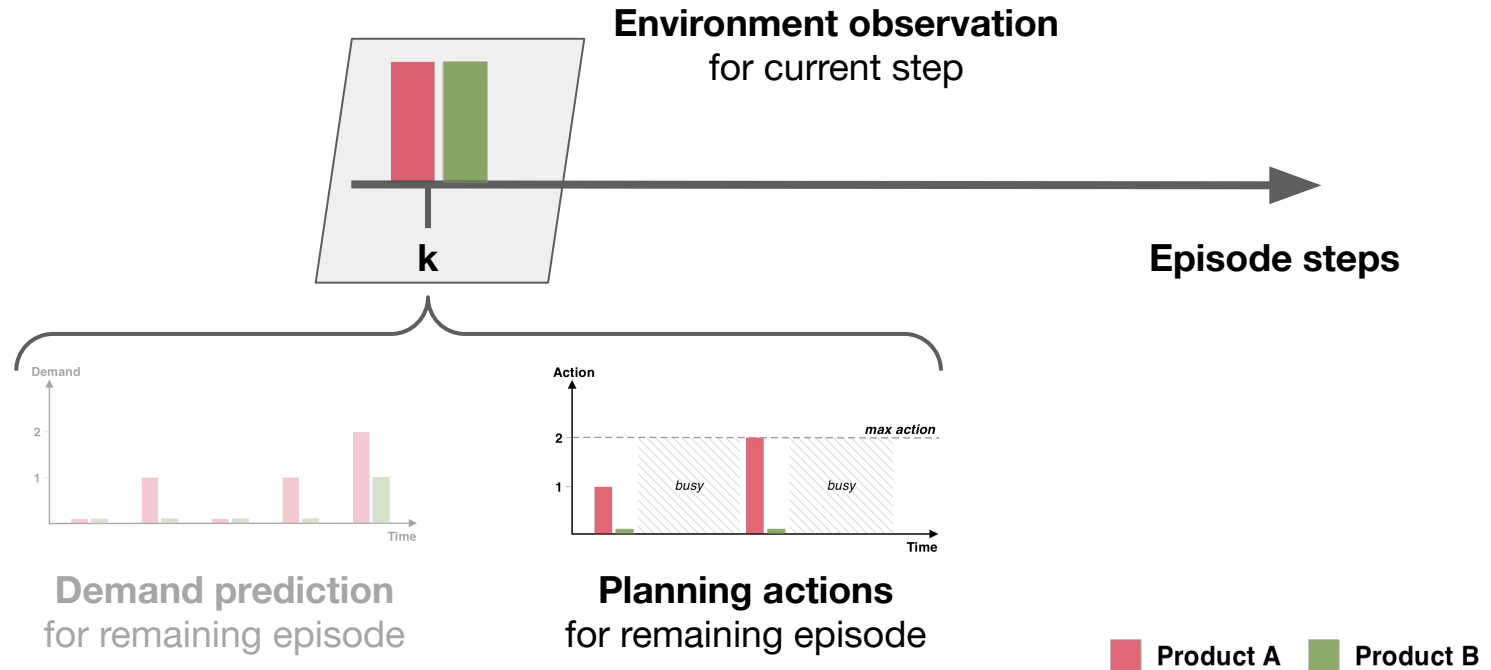
## Final model: Autoregressive Model (AR)

- predicts based on real consumer orders
- self-adjusts

# Demand Model - AR



# DP Agent



# Dynamic Programming

- **Discrete-time equation**

$$\tilde{x}_{k+1}^i = \tilde{x}_k^i + d_k^i - w_k^i$$

- Inventory  $\tilde{x}_k^i$
- Delivery  $d_k^i$
- Demand  $w_k^i$

- **Incorporated**

- Production time-shift
- Multi-products

# Dynamic Programming Algorithm

- **Discrete-time equation**

$$\tilde{x}_{k+1}^i = \tilde{x}_k^i + d_k^i - w_k^i$$

- Inventory  $\tilde{x}_k^i$
- Delivery  $d_k^i$
- Demand  $w_k^i$

- cost-to-go  $g_k$

- **Inventory Control**

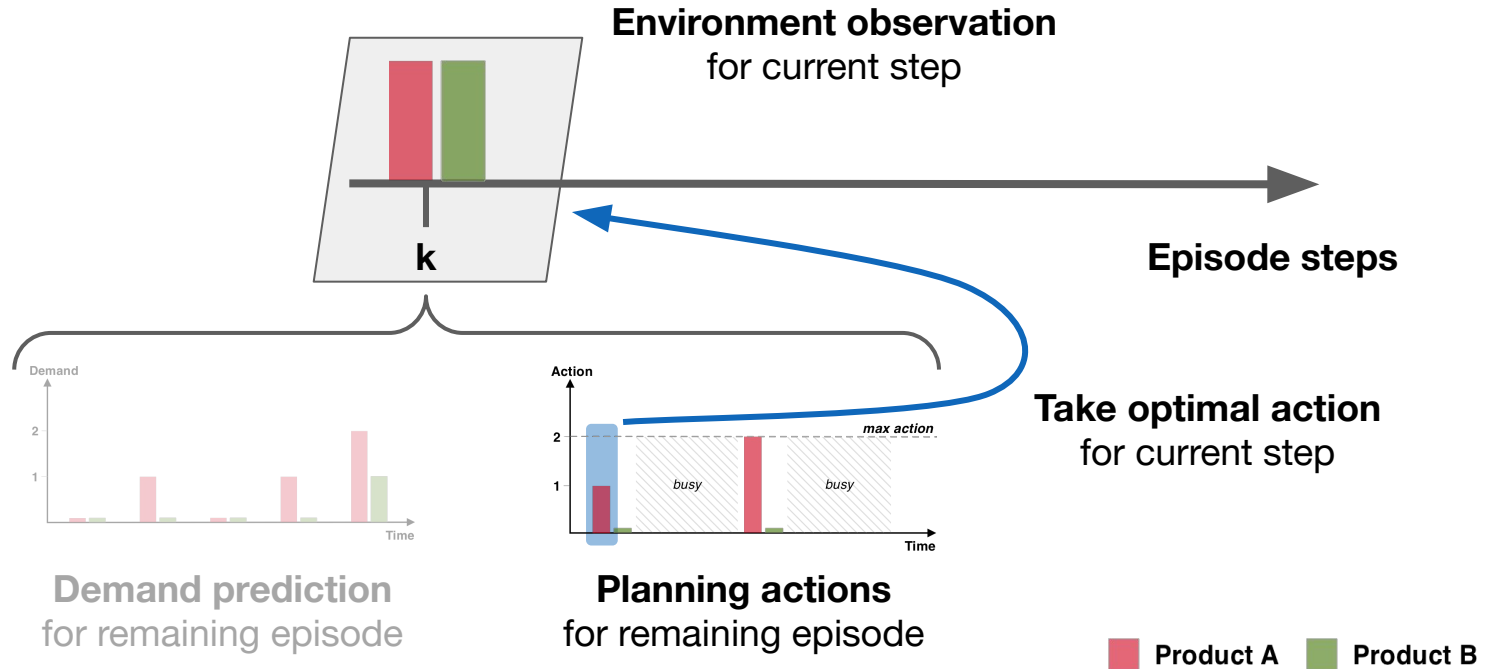
- Fixed demand to the prediction:  $\bar{w}_k^i$
- Solve deterministic problem:

$$\text{minimize } g_N(\tilde{x}_N^i) + \sum_k^{N-1} g_k(\tilde{x}_k^i, d_k^i, \bar{w}_k^i)$$

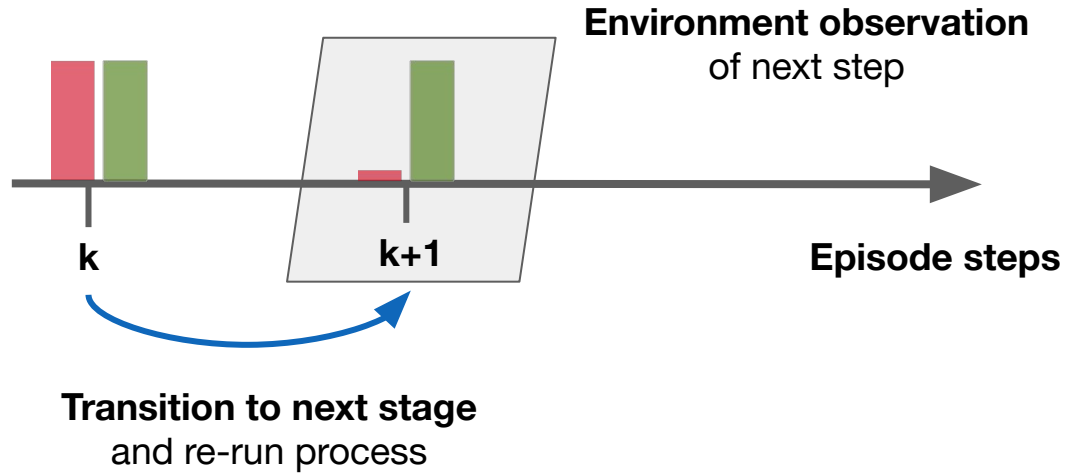
- **Uses *Principle of Optimality*** (Bellmann)

- Sequentially calculate optimal costs of tail subproblems
- Going from shorter to longer problems

# DP Agent



# DP Agent



■ Product A ■ Product B

# Results - Test cases

- In-distribution
  - Similar consumer demand (distribution) as training
- Out-of-distribution (domain transfer to new day / location)
  - **+50%** consumer demand **“Busy”**
  - **-50%** consumer demand **“Idle”**
  - **[-50%,+50%]** consumer demand **“Chaotic”**



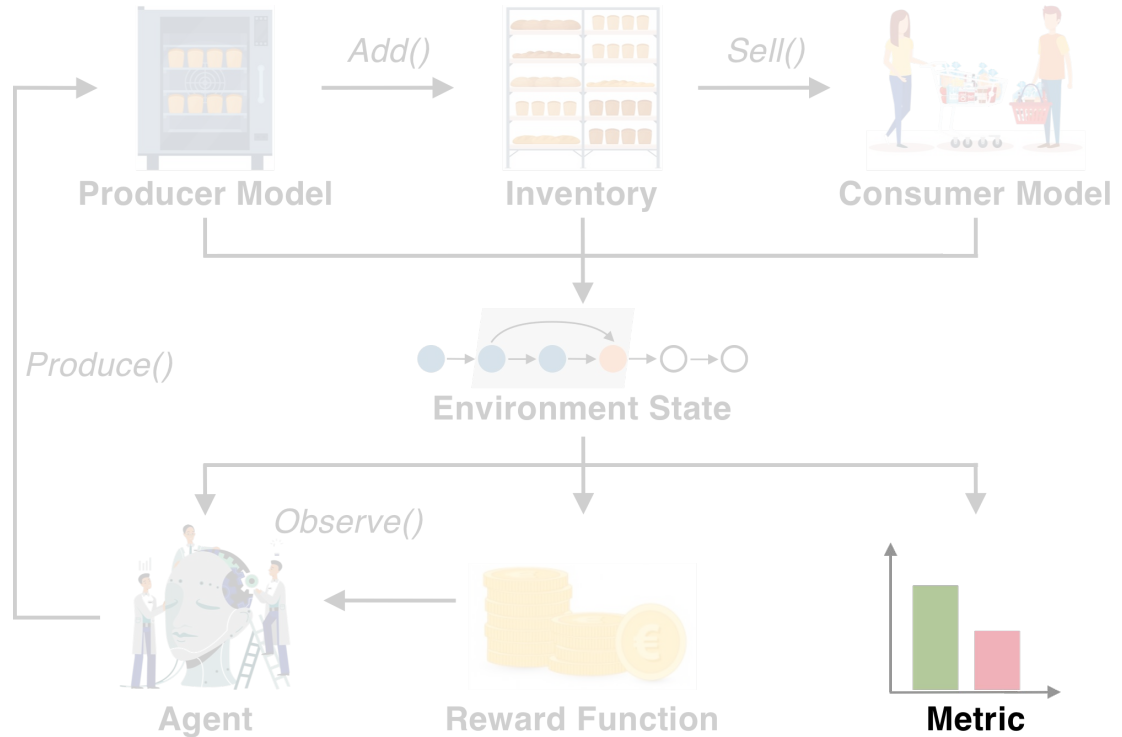
# Metrics

**Product sales:**  $S_{ratio}$

*sales out of total orders*

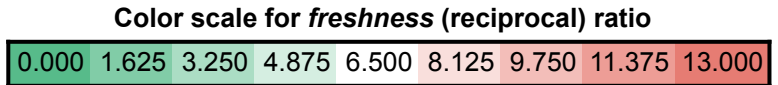
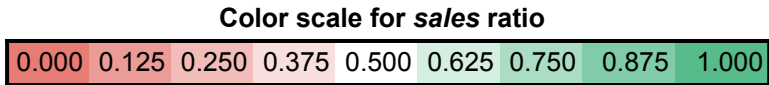
**Freshness:**  $p_{ratio}$   
(reciprocal)

*average age of the products*



# Results

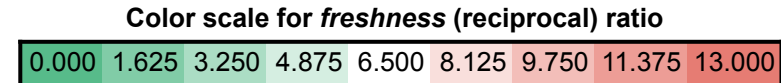
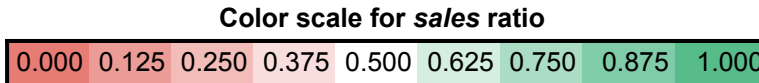
	<b>Sales</b>			<b>Freshness</b>			
	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
<b>Baseline</b>							
<b>DDPG</b>							
<b>DDPG w/ DR</b>							
<b>DP w/ AR</b>							
<b>DP w/ Oracle</b>							



# Results

	<b>Sales</b>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
<b>Baseline</b>	0.5503	0.4121	0.8188	0.5867
<b>DDPG</b>				
<b>DDPG w/ DR</b>				
<b>DP w/ AR</b>				
<b>DP w/ Oracle</b>				

	<b>Freshness</b>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
<b>Baseline</b>	5.3800	4.0300	10.1690	6.4300

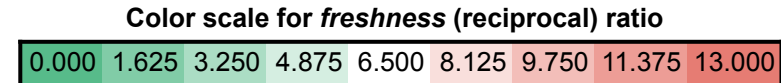
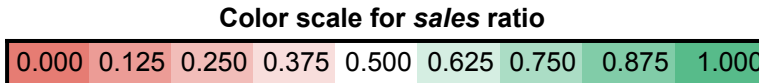


**Baseline** - Produces more than required, products less fresh

# Results

	<b>Sales</b>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	0.5503	0.4121	0.8188	0.5867
DDPG				
DDPG w/ DR				
DP w/ AR				
DP w/ Oracle	<b>0.6716</b>	<b>0.488</b>	<b>0.9521</b>	<b>0.6347</b>

	<b>Freshness</b>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	5.3800	4.0300	10.1690	6.4300
DDPG				
DDPG w/ DR				
DP w/ AR				
DP w/ Oracle	3.5250	2.7700	2.4500	3.6650



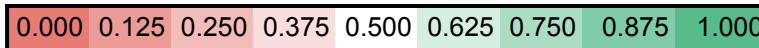
## DP with Oracle - Upper bound

# Results

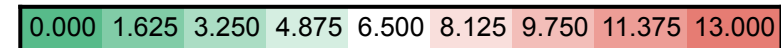
	<b>Sales</b>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	0.5503	0.4121	0.8188	0.5867
DDPG	0.4412	0.3316	0.7343	0.4222
DDPG w/ DR	0.2512	0.1902	0.4055	0.2597
DP w/ AR				
DP w/ Oracle	<b>0.6716</b>	<b>0.488</b>	<b>0.9521</b>	<b>0.6347</b>

	<b>Freshness</b>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	5.3800	4.0300	10.1690	6.4300
DDPG	2.5700	1.5399	6.5000	2.1500
DDPG w/ DR	<b>1.6750</b>	<b>0.9650</b>	<b>1.9200</b>	<b>1.2600</b>
DP w/ AR				
DP w/ Oracle	3.5250	2.7700	2.4500	3.6650

Color scale for sales ratio



Color scale for freshness (reciprocal) ratio



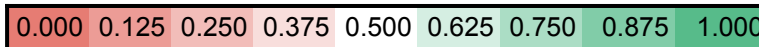
**DDPG, DDPG with Dynamics Randomization** - Produces less than required

# Results

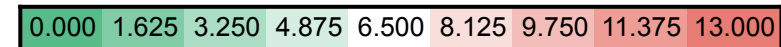
	<i>Sales</i>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	0.5503	0.4121	0.8188	0.5867
DDPG	0.4412	0.3316	0.7343	0.4222
DDPG w/ DR	0.2512	0.1902	0.4055	0.2597
DP w/ AR	0.5649	0.4013	0.8934	0.5736
DP w/ Oracle	<b>0.6716</b>	<b>0.488</b>	<b>0.9521</b>	<b>0.6347</b>

	<i>Freshness</i>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	5.3800	4.0300	10.1690	6.4300
DDPG	2.5700	1.5399	6.5000	2.1500
DDPG w/ DR	<b>1.6750</b>	<b>0.9650</b>	<b>1.9200</b>	<b>1.2600</b>
DP w/ AR	3.2736	1.6421	13.0000	3.6578
DP w/ Oracle	3.5250	2.7700	2.4500	3.6650

Color scale for sales ratio



Color scale for freshness (reciprocal) ratio



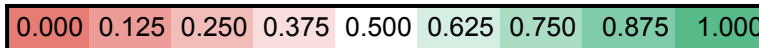
**DP with Autoregressive (AR) Prediction** - Balances well between both metrics

# Results

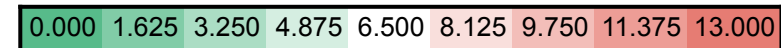
	<i>Sales</i>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	0.5503	0.4121	0.8188	0.5867
DDPG	0.4412	0.3316	0.7343	0.4222
DDPG w/ DR	0.2512	0.1902	0.4055	0.2597
DP w/ AR	0.5649	0.4013	0.8934	0.5736
DP w/ Oracle	<b>0.6716</b>	<b>0.488</b>	<b>0.9521</b>	<b>0.6347</b>

	<i>Freshness</i>			
	Similar demand	Higher mean demand (+50%)	Lower mean demand (-50%)	Random mean ([-50%,+50%])
Baseline	5.3800	4.0300	10.1690	6.4300
DDPG	2.5700	1.5399	6.5000	2.1500
DDPG w/ DR	<b>1.6750</b>	<b>0.9650</b>	<b>1.9200</b>	<b>1.2600</b>
DP w/ AR	3.2736	1.6421	13.0000	3.6578
DP w/ Oracle	3.5250	2.7700	2.4500	3.6650

Color scale for sales ratio



Color scale for freshness (reciprocal) ratio



**Overall** - good performance, improvements possible

# Summary

- **Brief Overview**

- Explore the available methods for domain transfer
- Create a consumer model suitable for the task
- Create a demand model that predicts the demands
- Implement Model-free (DDPG) and Model-based (DP) agents
- Compare agents with the baseline

- **Future Work**

- Improvement in consumer model with learning-based methods (e.g. LSTM)
- Approximate Dynamic Programming in case the state space is high



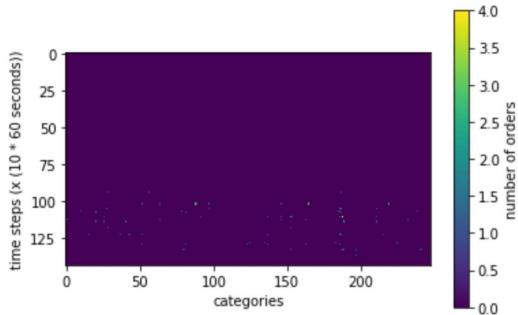
# Thank you for your attention! Questions?



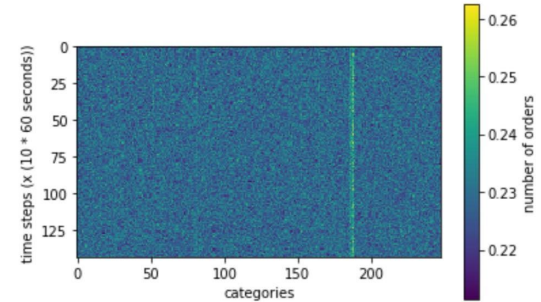
# Backup - VAE

## Take away Data

Training

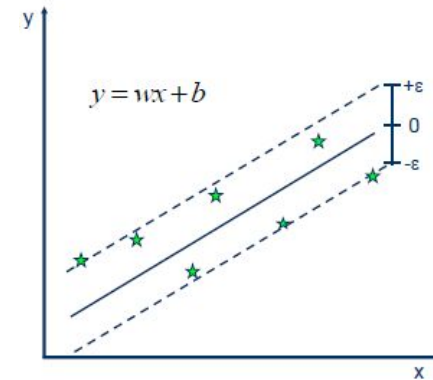


Generated



# Backup - Support Vector Regression

- Kernel Function maps data to higher dimensions
- At higher dimensions linear separation is possible
- A hyperplane between the data points
- A soft margin of tolerance for the hyperplane
- Minimizes the MSE (mean squared error)



• Solution:

$$\min \frac{1}{2} \|w\|^2$$

• Constraints:

$$y_i - wx_i - b \leq \varepsilon$$

$$wx_i + b - y_i \leq \varepsilon$$

# Backup - Autoregressive Model

$$X_t = c + \sum_{i=1}^p \varphi_i X_{t-i} + \epsilon_t$$

Here,  $\Phi$  are the parameters,  $X$  are the observations,  $\epsilon$  is white noise and  $c$  is the constant

- The output depends on the previous observations  $X$
- The parameters  $\Phi$  are adjusted
- Goal is to minimize MSE (mean squared error)

# Backup - Consumer models

- **Average**
  - For given time  $t$ , take average of all orders at this time
- **Random Nearest-Neighbor**
  - For given time  $t$ , go to nearest time  $t^*$ , and sample from all orders in  $t^*$  in a uniform random manner
- **Linear Regression / Support Vector Regression**
  - For  $n$  products, train  $n$  linear or support vector models

# Backup - Poisson Consumer Model

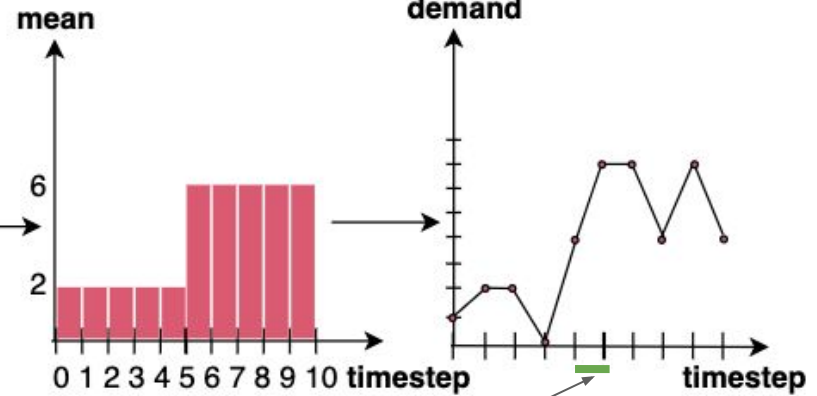
## Model Configuration

pizza:  
means: [10, 30]  
sandwich:  
means: [10, 10, 5, 5]  
...

10 , 30

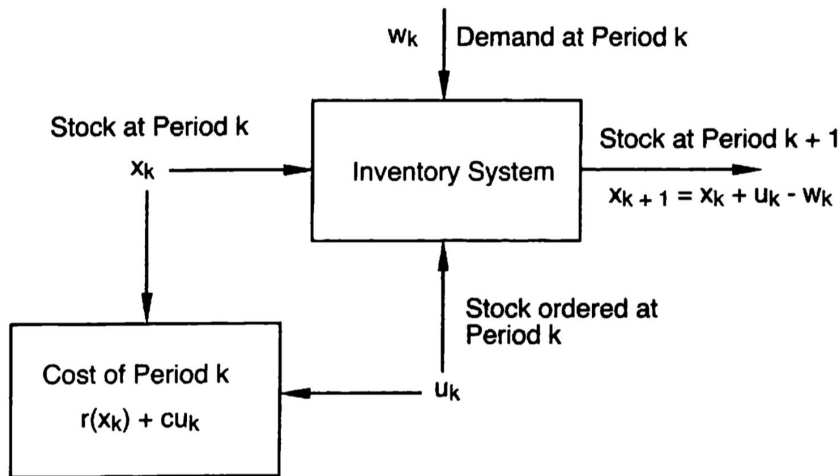
(number of time steps(N) = 10)

2, 2, 2, 2, 2, 6, 6, 6, 6, 6



$$\frac{(\lambda_i)^k e^{-\lambda_i}}{k!} = P(\text{k events in interval i})$$

# Backup - Original from D. Bertsekas



- Original example from D. Bertsekas

- Discrete-time system

$$x_{k+1} = f_k(x_k, u_k, w_k) = x_k + u_k - w_k$$

- Cost function that is additive over time

$$\begin{aligned} E \left\{ g_N(x_N) + \sum_{k=0}^{N-1} g_k(x_k, u_k, w_k) \right\} \\ = E \left\{ \sum_{k=0}^{N-1} (cu_k + r(x_k + u_k - w_k)) \right\} \end{aligned}$$

$$J_\pi(x_0) = E \left\{ R(x_N) + \sum_{k=0}^{N-1} (r(x_k) + c\mu_k(x_k)) \right\}$$

$\mu_k(x_k)$  = amount that should be ordered at time  $k$  if the stock is  $x_k$

# Backup - Metric

$$s_{ratio} = \frac{o_f}{o_f + o_u}$$

$o_f$  = number of fulfilled orders,

$o_u$  = number of unfulfilled orders

$$p_{ratio} = \frac{\sum_{i=0}^N a_i}{N}$$

$a_i$  = age of all products at time-step  $i$  in the inventory,

$N$  = number of time-steps.