

Planning and Control using Model-Based Reinforcement Learning

Pavel Czempin, Vincent Friedrich, Ruotong Liao, Jannik Nettelnstroth Mentors: Mathias Sundholm, Hamdi Belhassen Co-Mentor: Michael Rauchensteiner 24. February 2021



Introduction

Problem description

Challenge: Optimal production decisions in a bakery

Approach: Model-Based Reinforcement Learning

Goals:

- 1. Maximize sales
- 2. Minimize waste
- 3. Minimize shelf-time







ТШ

Our environment





Producer Model

- models an oven
- up to 30 products
- Agent can interact with this part

Inventory

• which products are ready for sale?

Consumer Model

- generates orders
- stochastic part of the environment



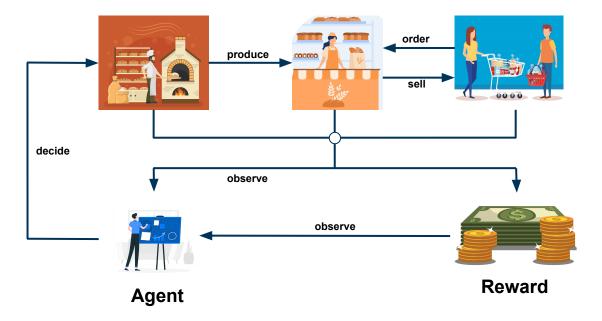
Complete setting

Agent

- decides when and which products to bake
- tries to maximize the reward

Reward

- measures performance of the agent with regard to the three goals
- is also used as metric to evaluate agent





PreciBake

Consumer models





Poisson model

→ samples orders from a poisson distribution with predefined mean



Real Data

→ customer data for 10 days in August and 15 days in December of the same year



Reward

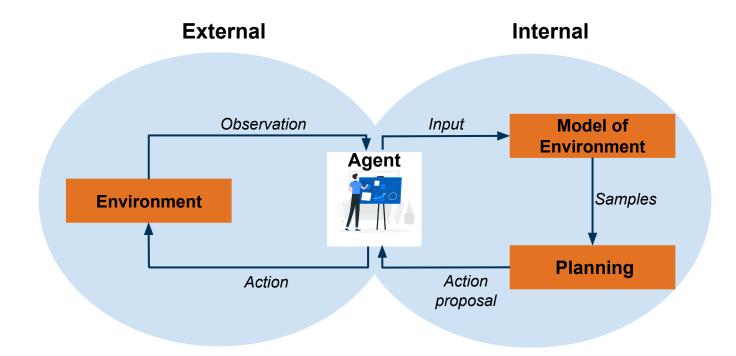


fulfilled-orders-ratio	$m_s = \frac{\# \text{sold products}}{\# \text{ordered products}}$
product-waste-ratio	$m_w = 1 - \frac{\#\text{wasted products}}{\#\text{produced products}}$
average_freshness	$m_f = \frac{1}{\# \text{sold products}} \sum_{\text{sold product } i} I[\text{age}_i \le \text{freshness_time}_i]$
overall reward	$m = \frac{4m_s + 4m_w + m_f}{9}$



Model-Based Reinforcement Learning







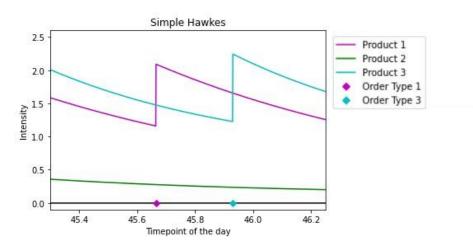
Hawkes Predictions

Simple Hawkes Model

- specified by Intensity function
 - Base rate μ
 - Excitement parameter α
 - Decay parameter ω
 - past orders *H(t)*

Intensity function
$$\lambda(t) = \mu + \alpha \sum_{q \in H(t)} e^{-\omega(t-q)}$$

• Orders of different products are independent



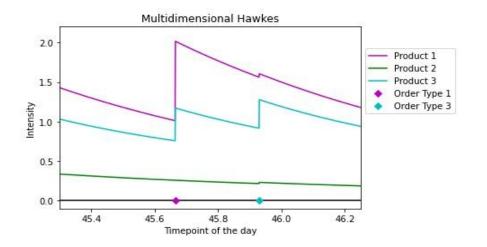


ΠΠ

Multidimensional Hawkes Model



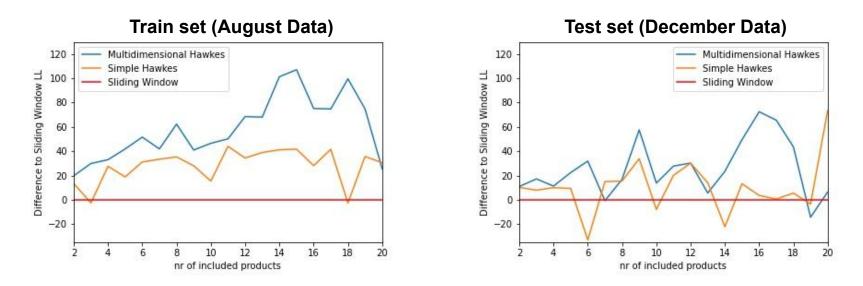
- **Extension** of the Simple Hawkes Model
- Intensity captures interdependencies between different products





Hawkes Experiments on Real Data



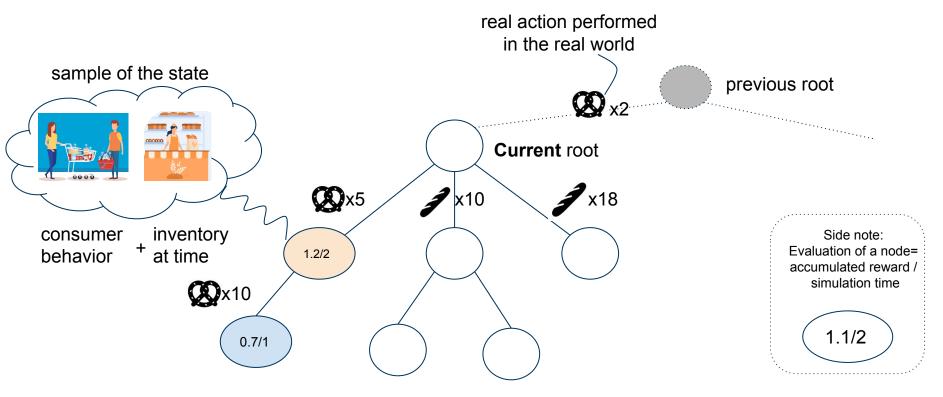


- Comparison of Multidimensional Hawkes, Simple Hawkes and Sliding Window (baseline model)
- Measure: relative average Log-Likelihood compared to baseline model



Monte Carlo Tree Search Planning



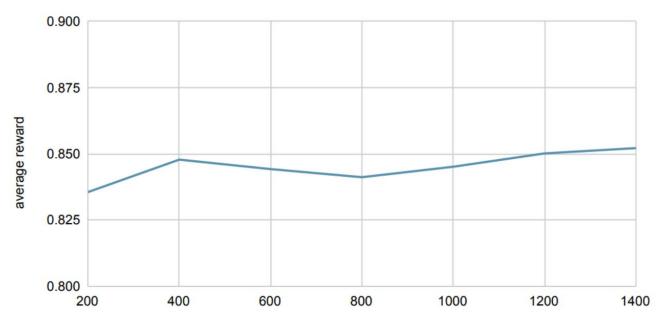






Heuristic idea: Simulation-based search algorithm relies on simulation budget

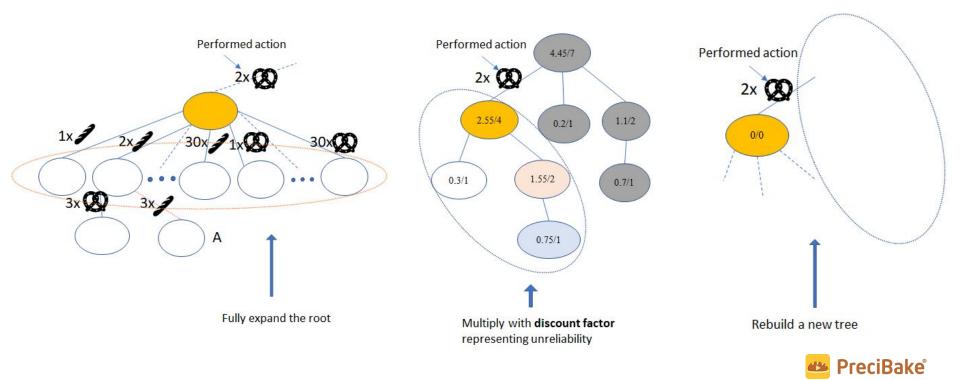
Hawkes MCTS Number of Simulations



PreciBake

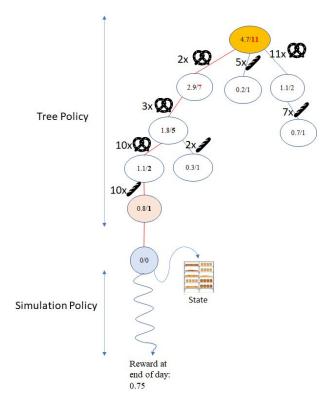


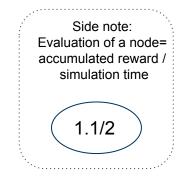
Tryouts to improve performance of the Monte Carlo Tree



ТЛ

Problem: The tree tends to explore deep leaf nodes and produces skewed tree

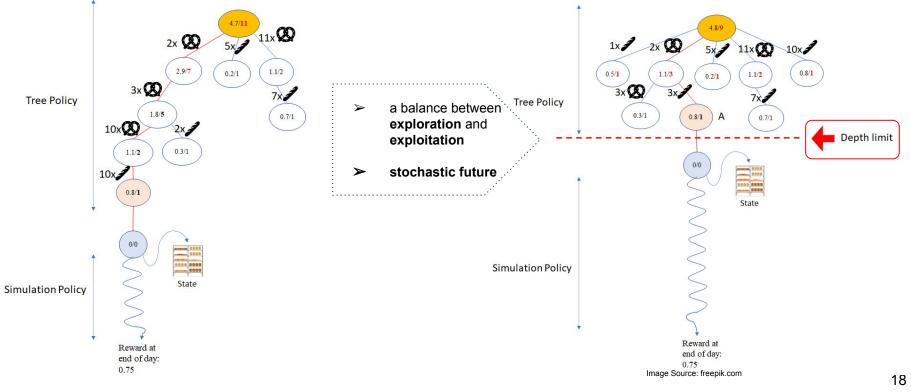






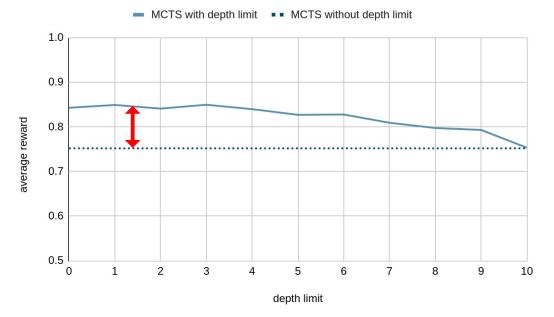


Problem:The tree tends to explore deep leaf nodes and produces skewed treeSolution:Introduce depth limitation to tree policy



ТШ

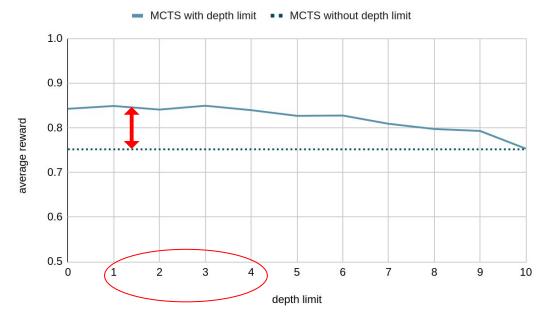
Hawkes MCTS Agent Depth Limit





ТШ

Hawkes MCTS Agent Depth Limit

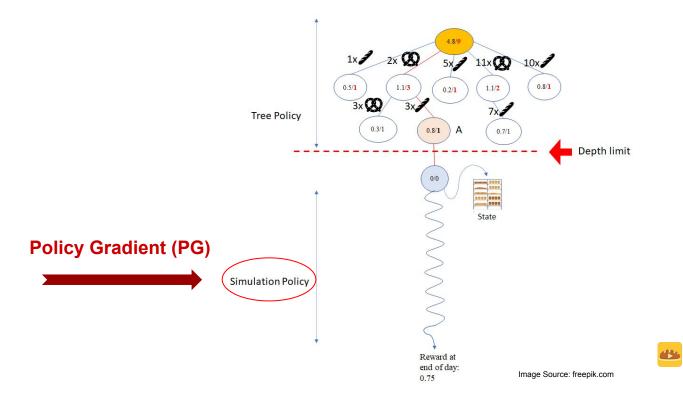






Problem: Random simulation policy might not provide accurate evaluation

Solution: Introduce RL simulation policy - Policy Gradients (PG)



PreciBake



Policy Gradient Method

- Policy function: $a \sim \pi(a|s)$
 - Sample actions for given state
- Directly optimize function approximator for likelihood of a sample
- Push up probabilities of actions that lead to higher return
- Proximal Policy Optimization (**PPO**)
 - Better convergence by clipping gradients



Final Experiments

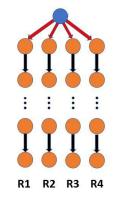




Threshold Agent



Monte Carlo Agent





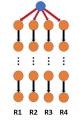
Experiment Setup - Agents



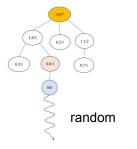
Threshold Agent



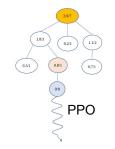
Monte Carlo Agent



Monte Carlo Tree Search Agent



MCTS+PG Agent





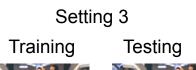
Experiment Setup - Environment











August

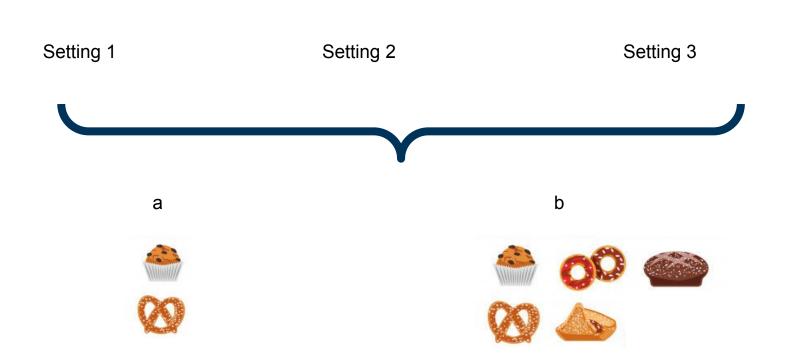


December





Image Source: freepik.com

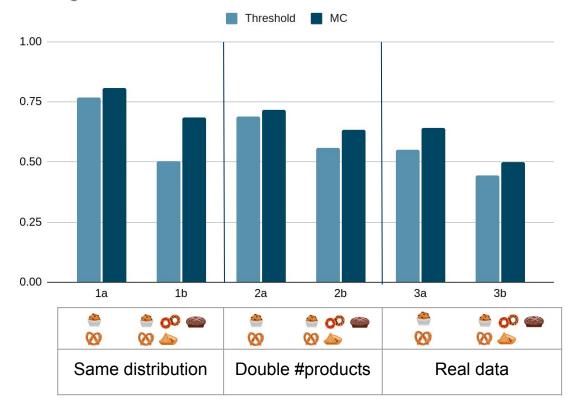


Experiment Setup





Average Rewards





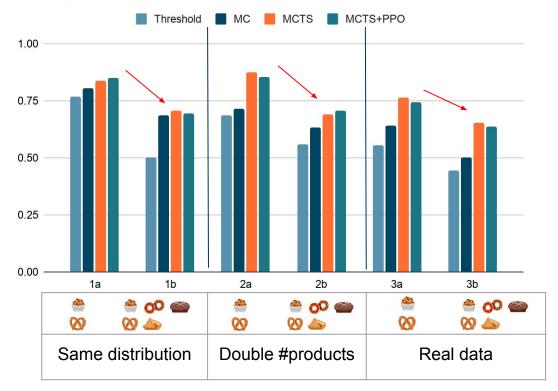
MCTS+PPO Threshold MC MCTS 1.00 0.75 0.50 0.25 0.00 1a 1b 2a 2b 3a 3b 🚔 🔿 🌚 🚔 📀 🌰 🐣 📀 🍩 00 8 00 0 00 00 / Double #products Same distribution Real data

Average Rewards



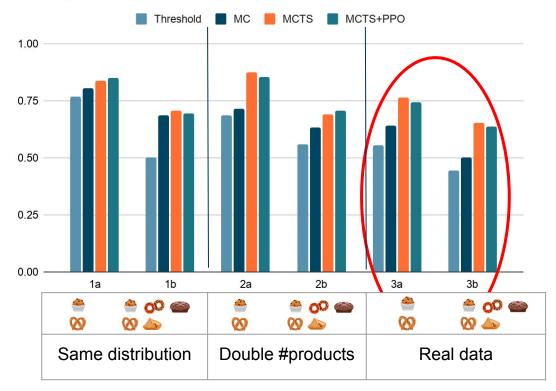


Average Rewards



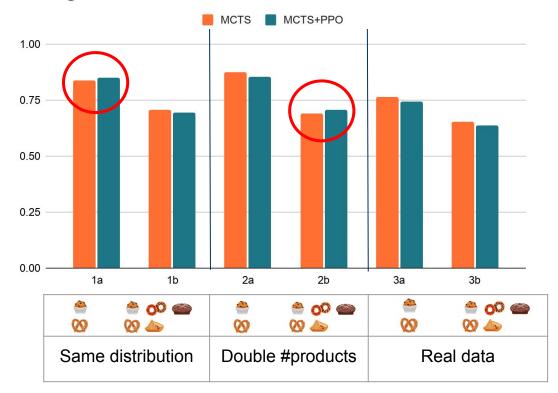


Average Rewards



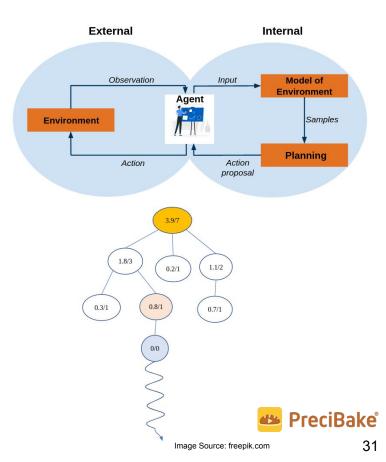


Average Rewards



Possible Future Improvements





- Improve internal model
- Expand usage of **RL** in MCTS
- Learn environment dynamics
- Model-based RL without tree search

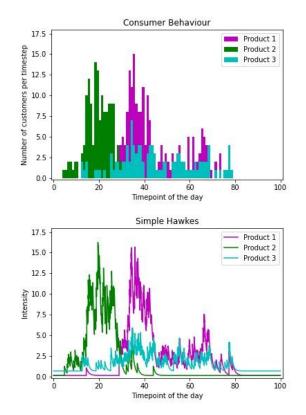
Thank you for your attention!

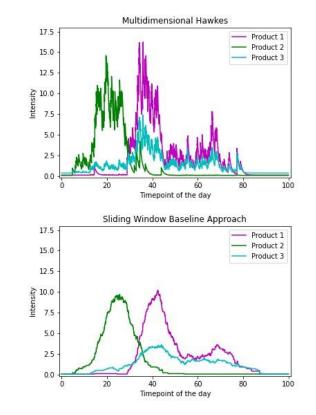
Backup slides

Hawkes Experiments on Real Data



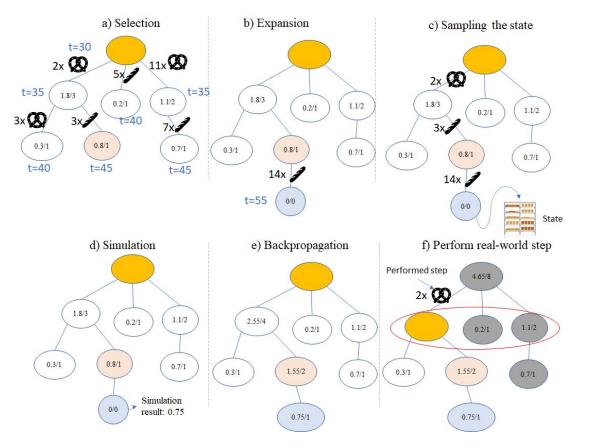
Real consumer behaviour and Intensity functions (example day, 3 products)















Sampling from Hawkes

Sampling via Thinning

Input: time t_i , history $\mathcal{H}(t_i)$ and end of episode t_{max} **Output:** tuple (t_{i+1}, j) of next order time t_{i+1} and product j or *None*

1. **Set** $t = t_i$

2. Choose
$$\mu_0 \in \mathbb{R}^+$$
 such that $\mu_0 \ge \sum_{k=1}^N \lambda_k (t' | \mathcal{H}(t_i)) \quad \forall t' \ge t$

```
3. Sample inter-event time \tau \sim \text{Exp}(1/\mu_0)
```

4. **Set** $t = t + \tau$

5. If
$$t > t_{\max}$$
:
return None

6. Sample
$$j \sim \text{Categorical}\left(\frac{\lambda_1(t|\mathcal{H}(t_i))}{\mu_0}, ..., \frac{\lambda_N(t|\mathcal{H}(t_i))}{\mu_0}, 1 - \frac{\sum_{k=1}^N \lambda_k(t|\mathcal{H}(t_i))}{\mu_0}\right)$$

If $j \in \{1, ..., N\}$:
return (t, j)
Else:
Go to 3.

